# Lecture 18: Autoencoders
# CS480/680 Intro to Machine Learning

Pascal Poupart
David R. Cheriton School of Computer Science
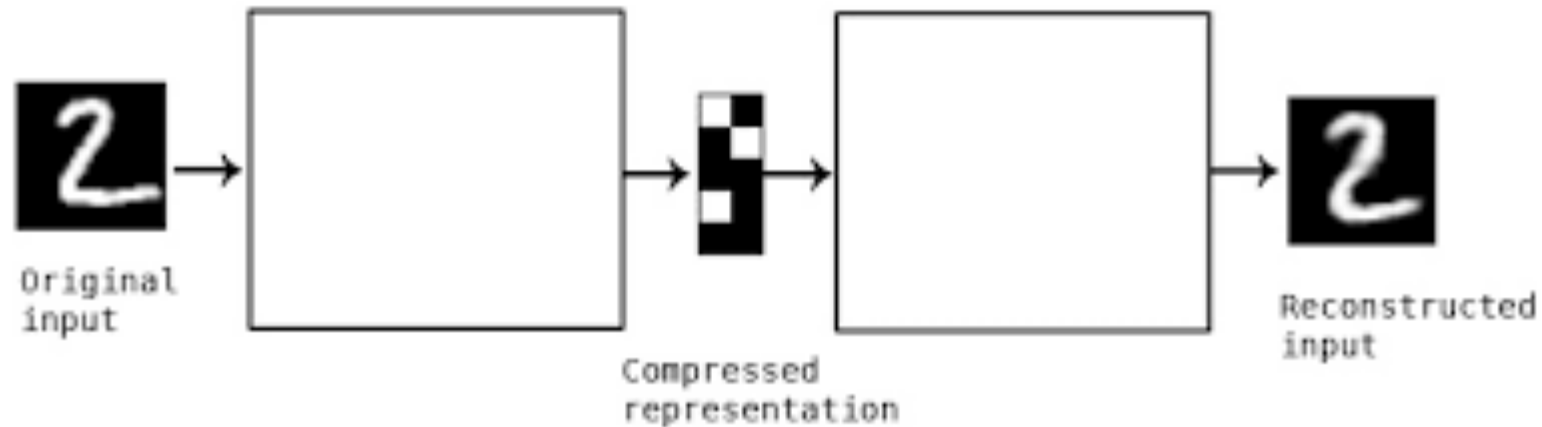
UNIVERSITY OF
WATERLOO

# Autoencoder

- Special type of feed forward network for
    - Compression
    - Denoising
    - Sparse representation
    - Data generation

# Autoencoder

- Encoder: $f(\quad)$

- Decoder: $g(\quad)$

- Autoencoder: $g\big(f(\boldsymbol{x})\big) = \boldsymbol{x}$



Original input     Compressed representation     Reconstructed input

# Linear Autoencoder

- $f$ and $g$ are linear

  - Matrix representations: $\boldsymbol{W}_f$ and $\boldsymbol{W}_g$

- Picture:

# Linear Autoencoder

- Objective: find weights $\boldsymbol{W}_f$ and $\boldsymbol{W}_g$ that minimize reconstruction error

$$\min_{\boldsymbol{W}} \frac{1}{2} \sum_n \left\| \boldsymbol{W}_g \boldsymbol{W}_f \boldsymbol{x}_n - \boldsymbol{x}_n \right\|_2^2$$

- Algorithm: backpropagation
  - Gradient descent

- When using Euclidean norm (i.e., squared loss), solution is the same as principal component analysis (PCA)

UNIVERSITY OF
WATERLOO

# Principal Component Analysis

- Hidden nodes: compressed representation

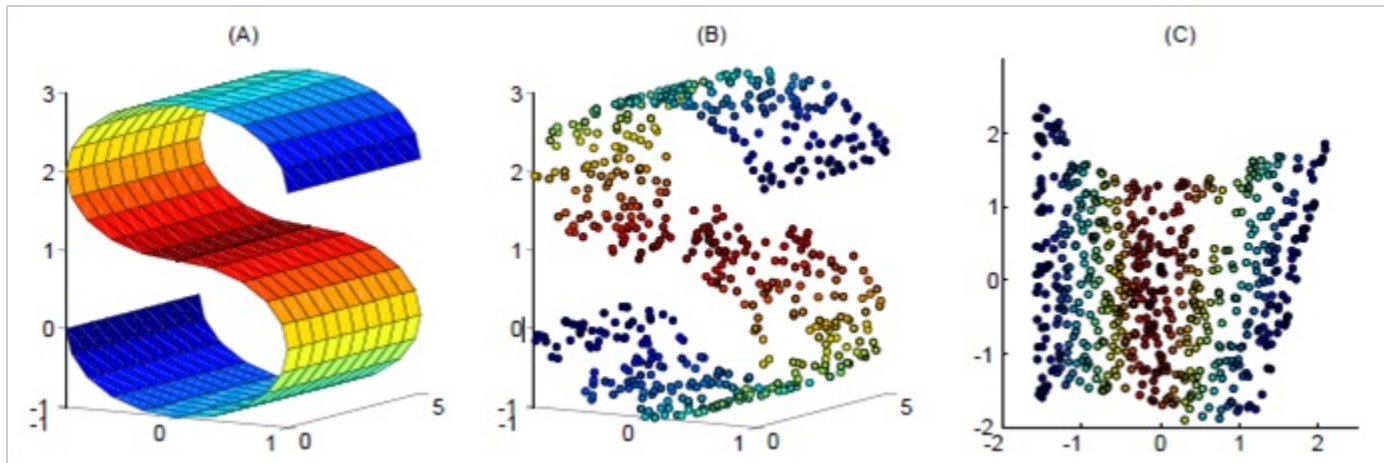# Nonlinear Autoencoder

- $f$ and $g$ are non-linear functions

$$\min_W \frac{1}{2} \sum_n \left\| g(f(x_n; W_f); W_g) - x_n \right\|_2^2$$

- Hidden nodes: non-linear manifold

UNIVERSITY OF
WATERLOO

# Deep Autoencoders

- $f$ and $g$ often consist of multiple layers

- In theory, one hidden layer in $f$ and $g$ is sufficient to represent any possible compression

- Multiple hidden layers in $f$ and $g$ is often better

UNIVERSITY OF
WATERLOO

# Sparse Representations

- When more hidden nodes than inputs, use regularization to constrain autoencoder

- Example: force hidden nodes to be sparse

$$\min_{\boldsymbol{W}} \frac{1}{2} \sum_n \left|\left| g(f(\boldsymbol{x}_n; \boldsymbol{W}_f); \boldsymbol{W}_g) - \boldsymbol{x}_n \right|\right|_2^2 + c \, nnz\left(f(\boldsymbol{x}_n; \boldsymbol{W}_f)\right)$$

**Sparse hidden nodes**

where $nnz\left(f(\boldsymbol{x}_n; \boldsymbol{W}_f)\right)$ is the number of non-zero entries in the vector produced by $f$.

- Approximate objective: L1 regularization

$$\min_{\boldsymbol{W}} \frac{1}{2} \sum_n \left|\left| g(f(\boldsymbol{x}_n; \boldsymbol{W}_f); \boldsymbol{W}_g) - \boldsymbol{x}_n \right|\right|_2^2 + c \, \left|\left| f(\boldsymbol{x}_n; \boldsymbol{W}_f) \right|\right|_1$$

UNIVERSITY OF
WATERLOO

# Denoising Autoencoder

- Consider noisy version $\tilde{x}$ of the input $x$

- Data denoising

$$\min_{W} \frac{1}{2} \sum_{n} \left|\left| g\left(f(\tilde{x}_n; W_f); W_g\right) - x_n \right|\right|_2^2 + c \left|\left| f(\tilde{x}_n; W_f) \right|\right|_1$$



$$x \qquad\qquad \tilde{x} \qquad\qquad g(f(\tilde{x}))$$

original      perturbed      reconstructed

# Probabilistic Autoencoder

- Let $f$ and $g$ represent conditional distributions

$$f: \Pr(\boldsymbol{h}|\boldsymbol{x}; \boldsymbol{W}_f) \quad \text{and} \quad g: \Pr(\boldsymbol{x}|\boldsymbol{h}; \boldsymbol{W}_g)$$

   by using sigmoid, softmax or linear units at the hidden and output layers

- Picture

UNIVERSITY OF
**WATERLOO**

# Generative Model

- Sample $h$ from some distribution $\Pr(h)$

- Sample $x$ from decoder $\Pr(x|h; W_g)$



$\Pr(h)$

$\Pr(x|h; W_g)$

Noise ~ N(0,1)

Generative Model

UNIVERSITY OF
WATERLOO