

Lecture 18: Autoencoders

CS480/680 Intro to Machine Learning

2023-3-16

Pascal Poupart
David R. Cheriton School of Computer Science

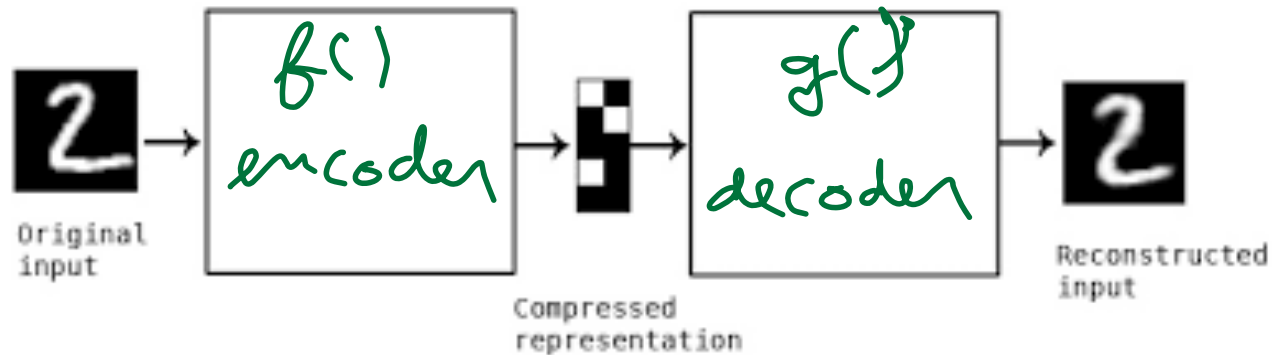


Autoencoder

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

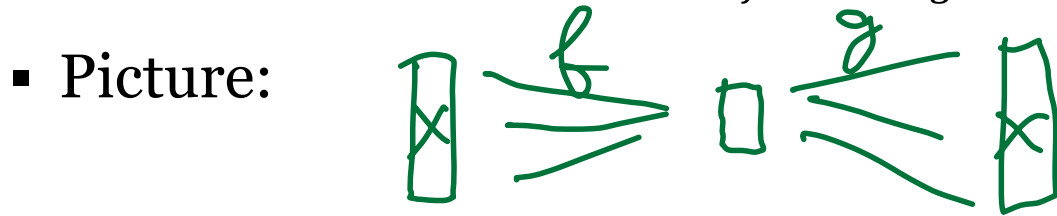
Autoencoder

- Encoder: $f(\)$
- Decoder: $g(\)$
- Autoencoder: $g(f(x)) = x$



Linear Autoencoder

- f and g are linear
 - Matrix representations: W_f and W_g



$$\boxed{x^T} \quad \boxed{W_f^T} \quad \boxed{W_g^T} = \boxed{\hat{x}^T}$$

Linear Autoencoder

- Objective: find weights W_f and W_g that minimize reconstruction error

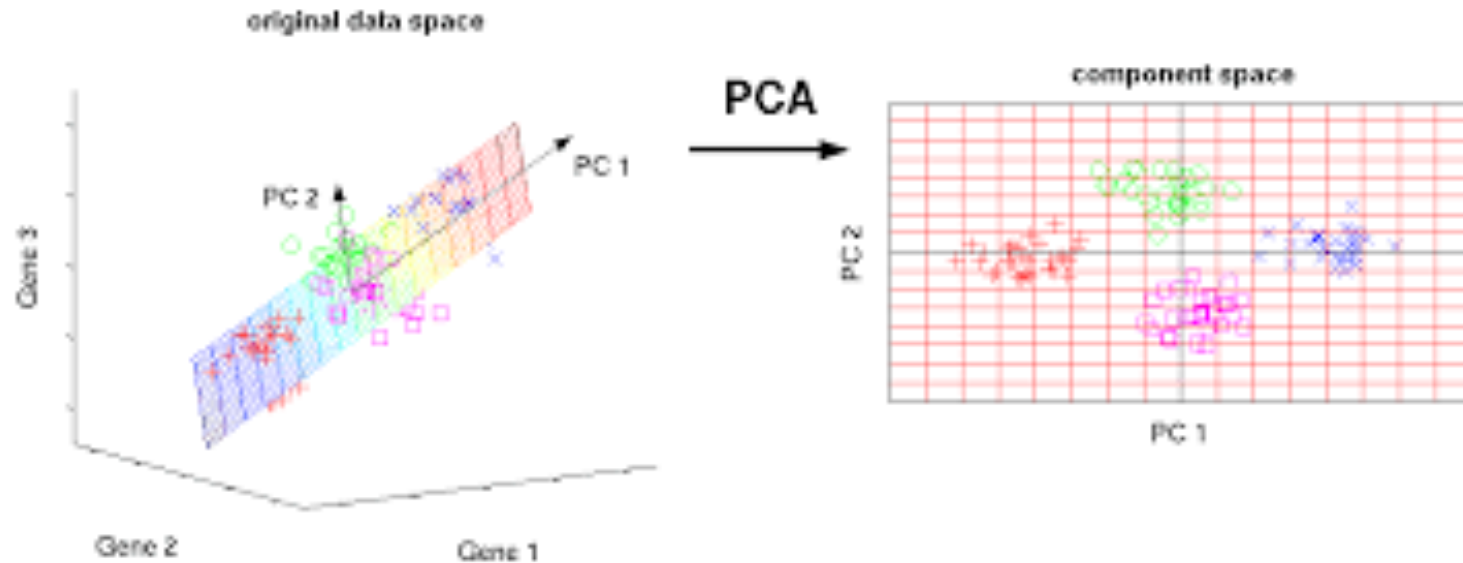
$$\min_W \frac{1}{2} \sum_n \|W_g W_f x_n - x_n\|_2^2$$

- Algorithm: backpropagation
 - Gradient descent

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

Principal Component Analysis

- Hidden nodes: compressed representation

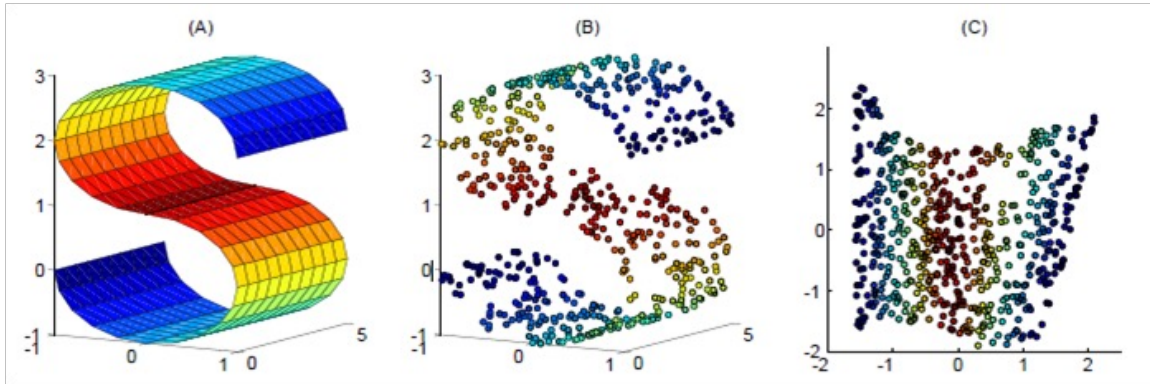


Nonlinear Autoencoder

- f and g are non-linear functions

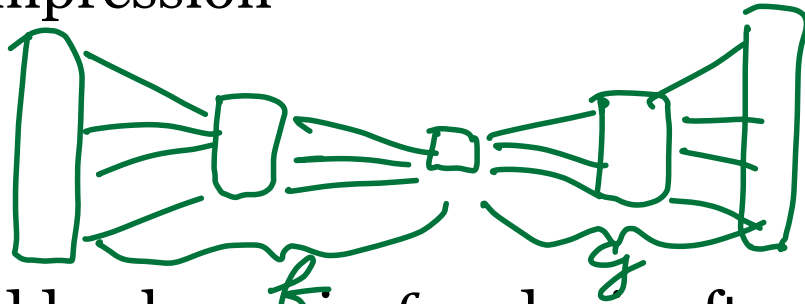
$$\min_W \frac{1}{2} \sum_n \left\| g(f(\mathbf{x}_n; \mathbf{W}_f); \mathbf{W}_g) - \mathbf{x}_n \right\|_2^2$$

- Hidden nodes: non-linear manifold

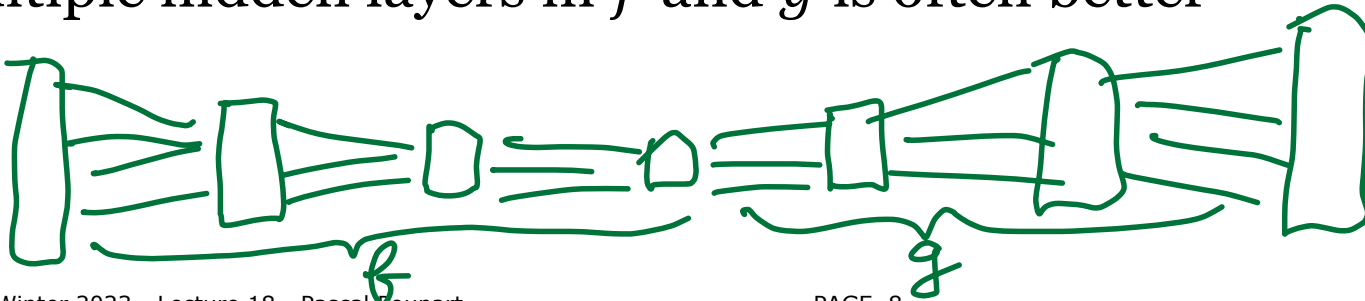


Deep Autoencoders

- 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



Sparse Representations

- When more hidden nodes than inputs, use regularization to constrain autoencoder
- Example: force hidden nodes to be sparse

$$\min_{\mathbf{W}} \frac{1}{2} \sum_n \left\| g(f(\mathbf{x}_n; \mathbf{W}_f); \mathbf{W}_g) - \mathbf{x}_n \right\|_2^2 + c \underbrace{nnz(f(\mathbf{x}_n; \mathbf{W}_f))}_{\text{Sparse hidden nodes}}$$

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

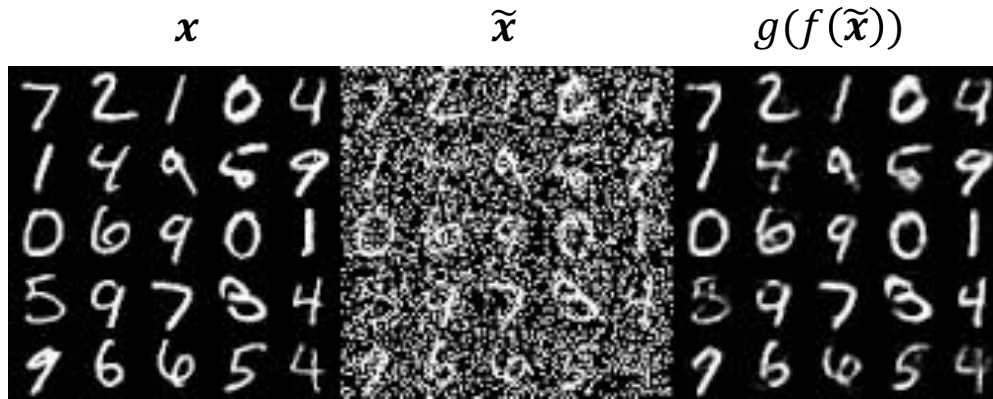
- Approximate objective: L1 regularization

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

Denoising Autoencoder

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

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



original

perturbed

reconstructed

Probabilistic Autoencoder

- Let f and g represent conditional distributions

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

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

- Picture



Sigmoid: Bernoulli $P(h): h \in \{0, 1\}$
Softmax: Categorical $P(h): h \in \text{one hot vectors}$
Linear: Gaussian $P(h|\mu, \sigma^2)$ → variance
mean ←

Generative Model

- Sample \mathbf{h} from some distribution $\Pr(\mathbf{h})$
- Sample \mathbf{x} from decoder $\Pr(\mathbf{x}|\mathbf{h}; \mathbf{W}_g)$

