# Autoencoders

Lecture slides for Chapter 14 of *Deep Learning*

www.deeplearningbook.org

Ian Goodfellow

2016-09-30

# Structure of an Autoencoder

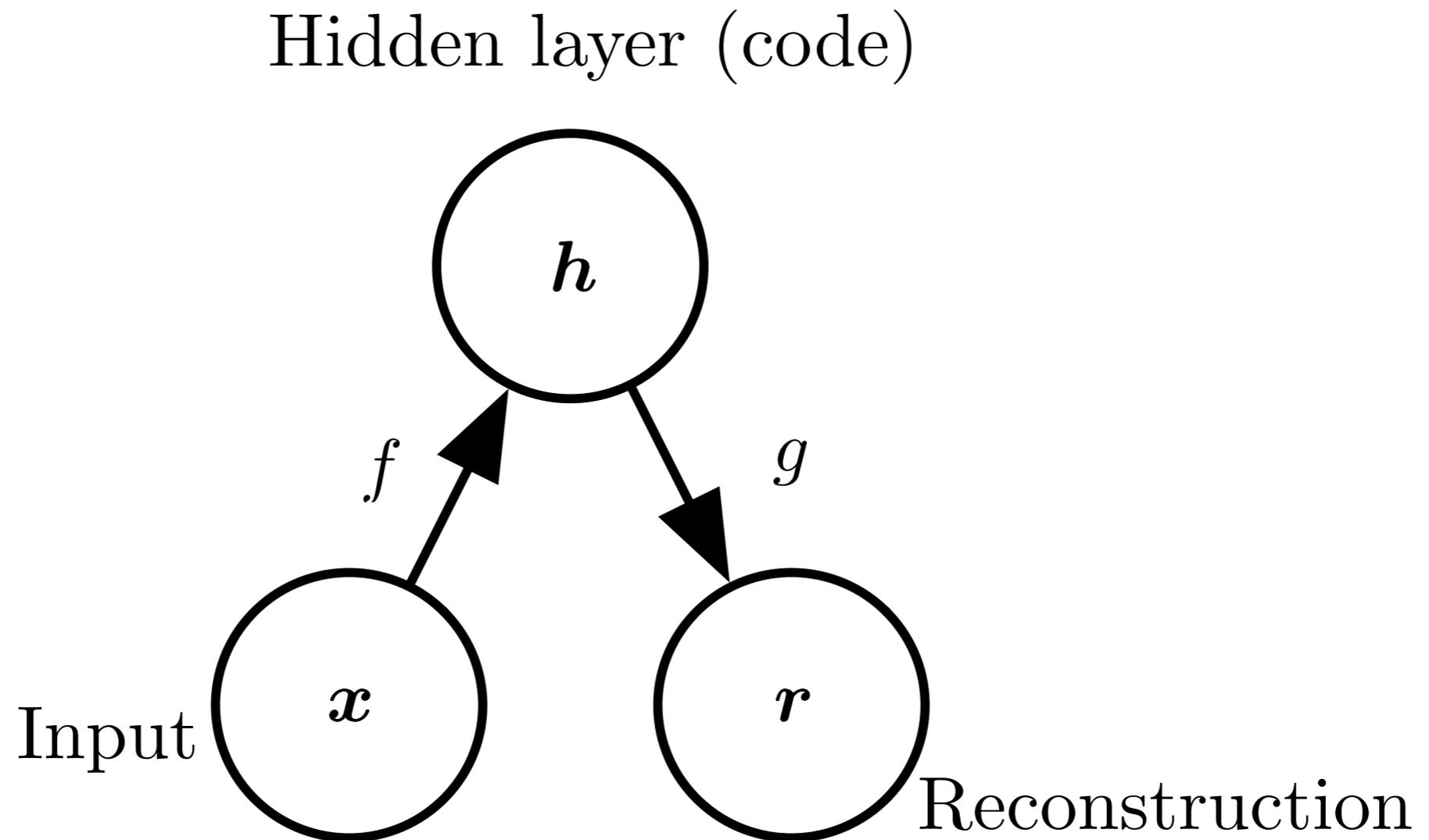Hidden layer (code)



Figure 14.1
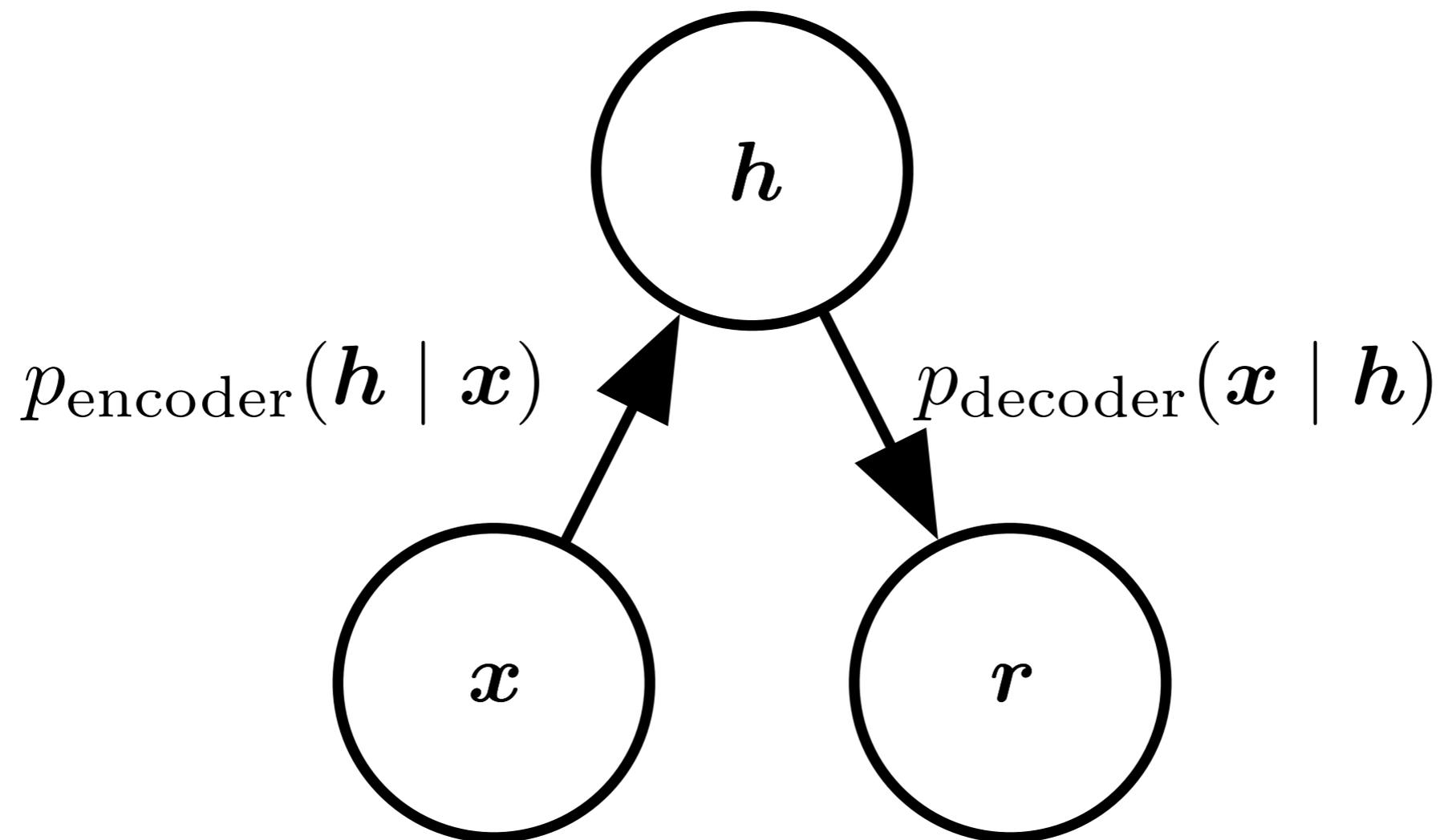
# Stochastic Autoencoders



Figure 14.2

# Avoiding Trivial Identity

- Undercomplete autoencoders

  - $h$ has lower dimension than $x$

  - $f$ or $g$ has low capacity (e.g., linear $g$)

  - Must discard some information in $h$

- Overcomplete autoencoders

  - $h$ has higher dimension than $x$

  - Must be regularized

# Regularized Autoencoders

- Sparse autoencoders

- Denoising autoencoders

- Autoencoders with dropout on the hidden layer

- Contractive autoencoders

# Sparse Autoencoders

- Limit capacity of autoencoder by adding a term to the cost function penalizing the code for being larger

- Special case of variational autoencoder

  - Probabilistic model

  - Laplace prior corresponds to L1 sparsity penalty
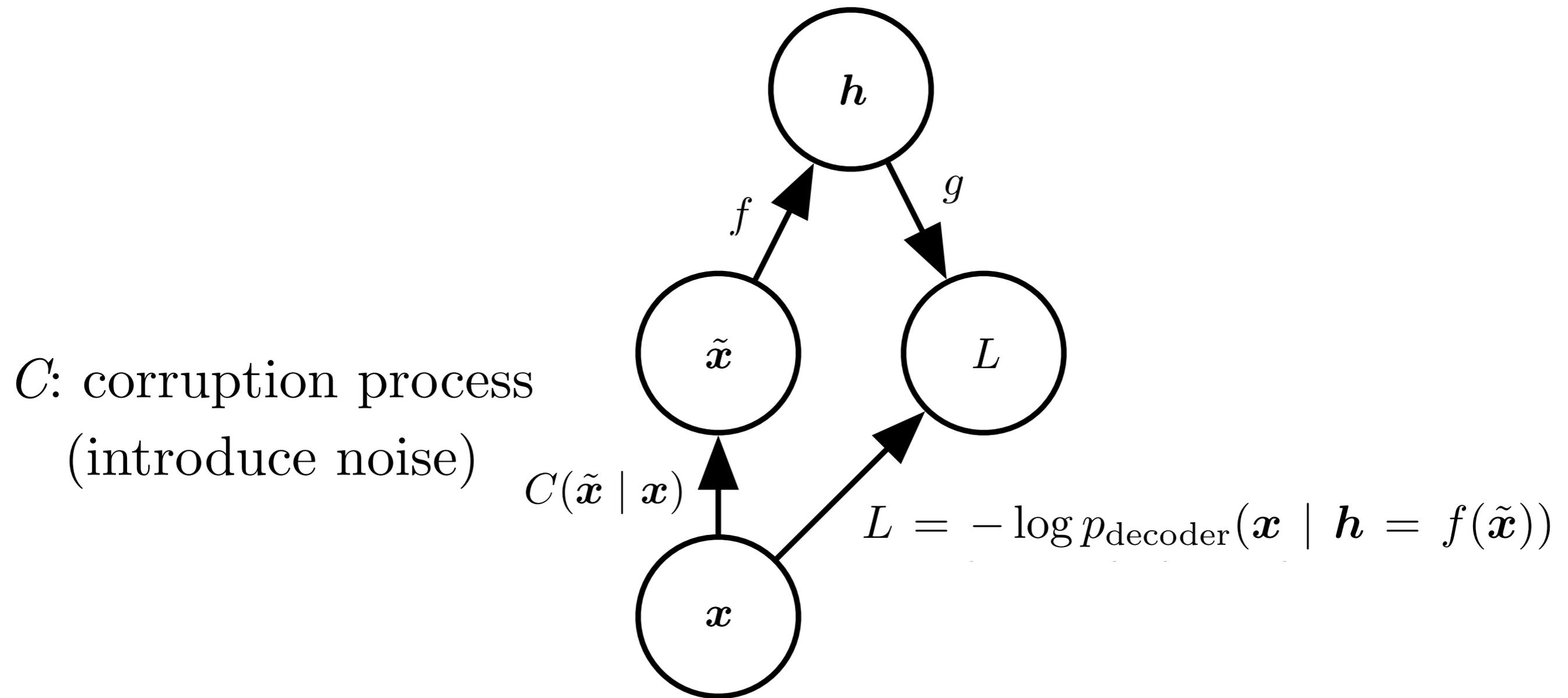
  - Dirac variational posterior

# Denoi~~s~~



$C$: corruption proce~~ss~~
   (introduce noise)

$$L = -\log p_{\text{decoder}}(\boldsymbol{x} \mid \boldsymbol{h} = f(\tilde{\boldsymbol{x}}))$$

# Denoising Autoencoders Learn a Manifold



Figure 14.4

# Score Matching

- Score: $\nabla_{\boldsymbol{x}} \log p(\boldsymbol{x})$. $\hspace{4cm}$ (14.15)

- Fit a density model by matching score of model to score of data

- Some denoising autoencoders are equivalent to score matching applied to some density models

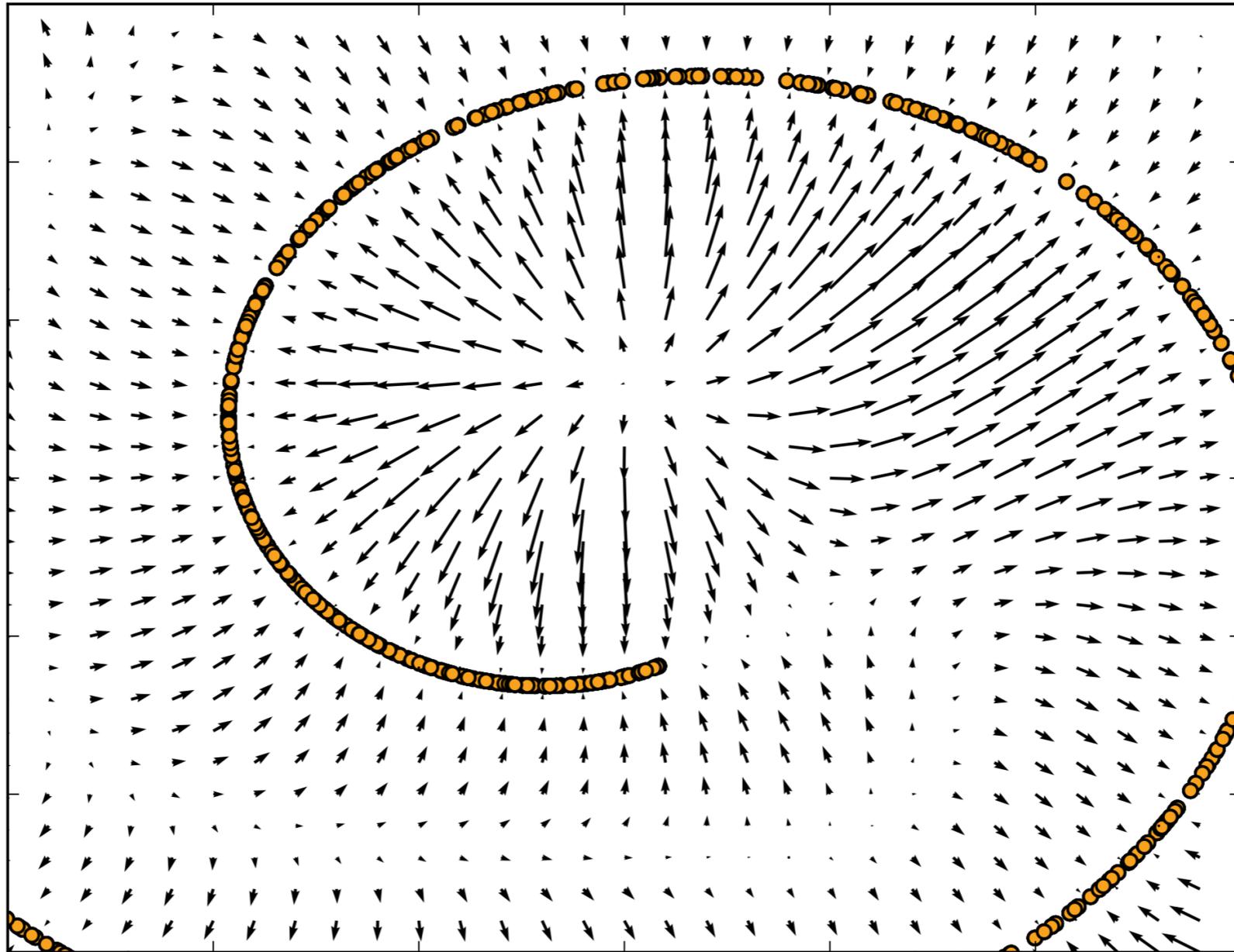# Vector Field Learned by a Denoising Autoencoder



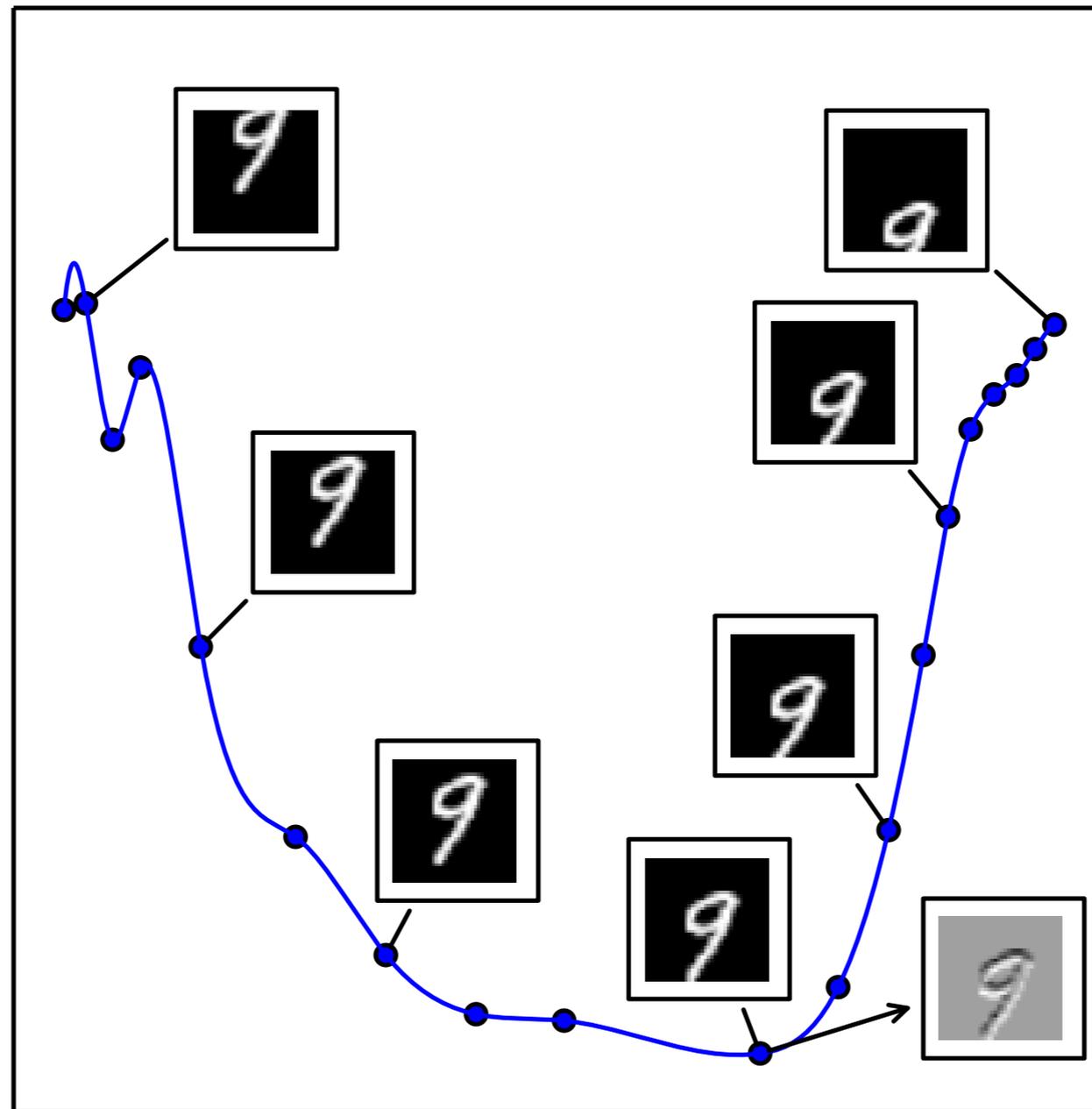Figure 14.5

# Tangent Hyperplane of a Manifold



Figure 14.6

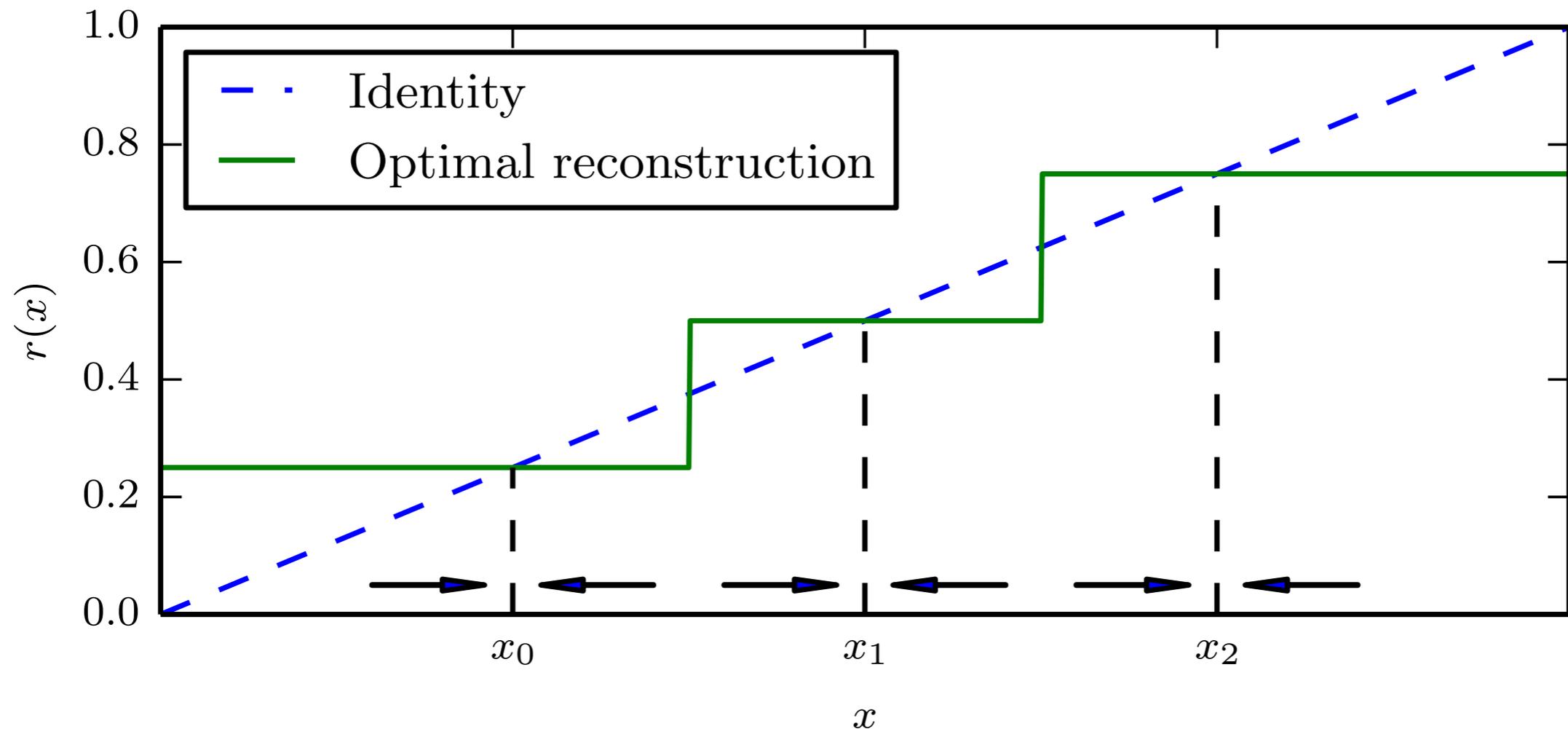# Learning a Collection of 0-D Manifolds by Resisting Perturbation



Figure 14.7

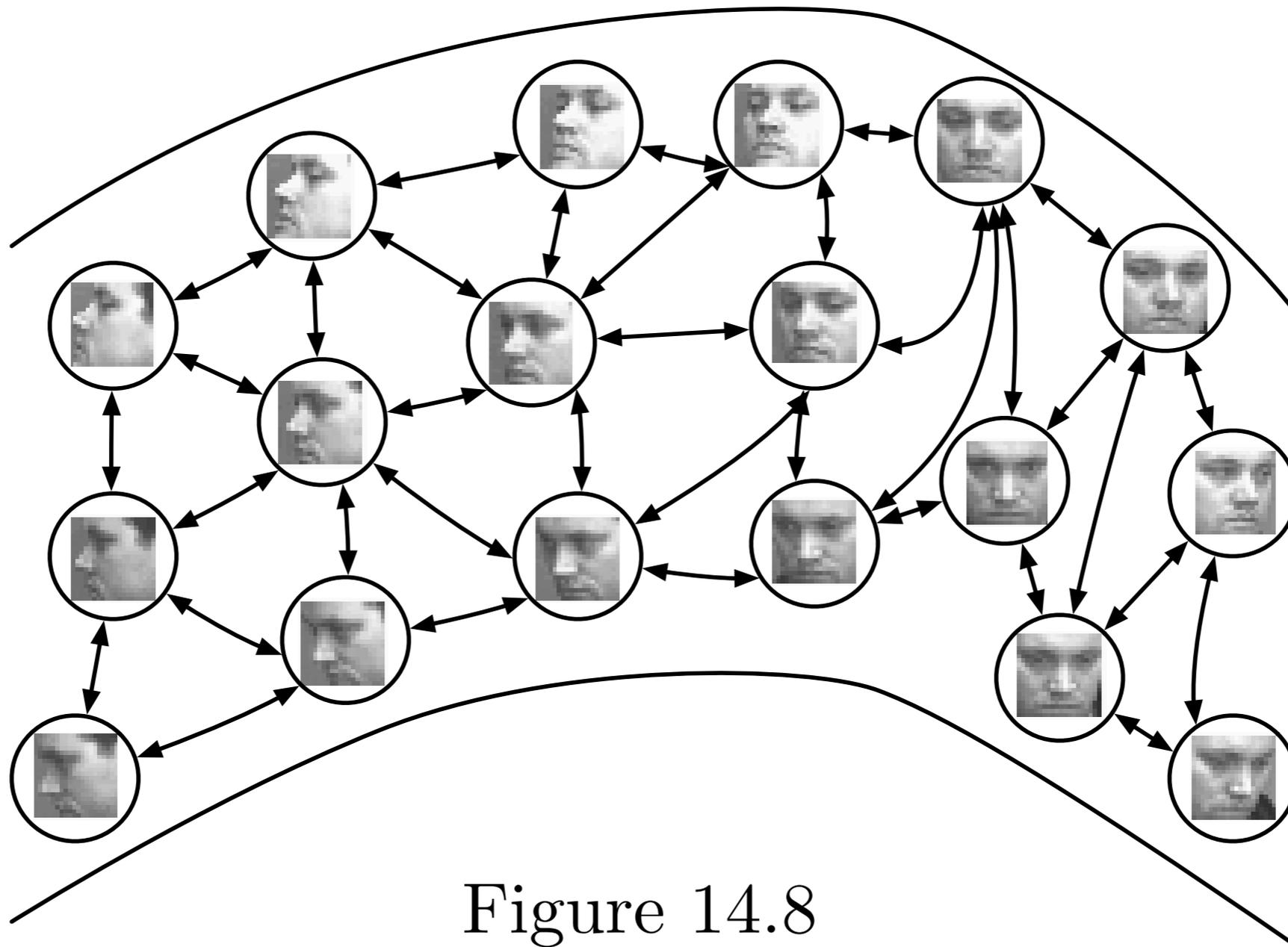# Non-Parametric Manifold Learning with Nearest-Neighbor Graphs



Figure 14.8

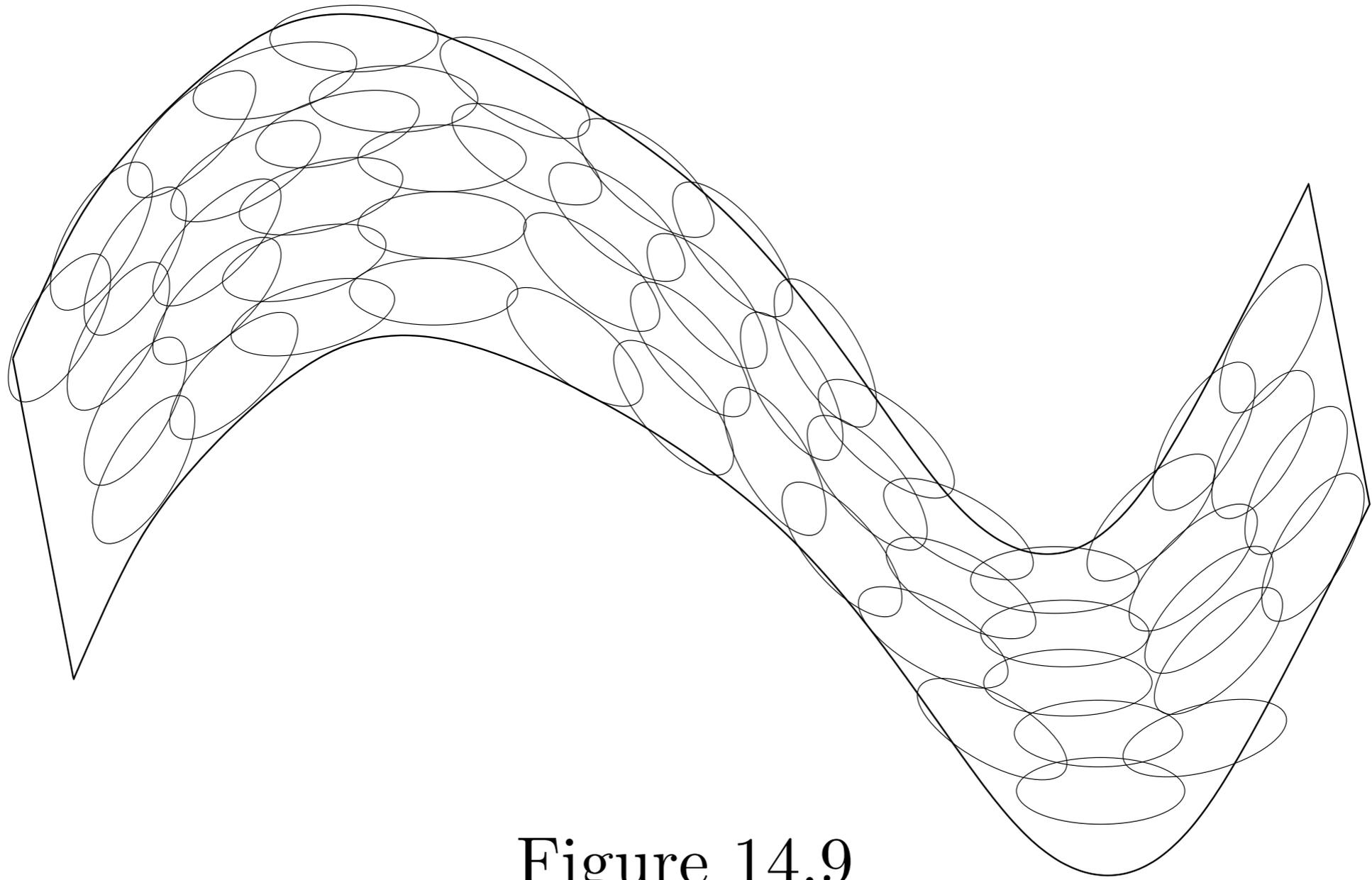# Tiling a Manifold with Local Coordinate Systems



Figure 14.9

# Contractive Autoencoders

$$\Omega(\boldsymbol{h}) = \lambda \left\| \frac{\partial f(\boldsymbol{x})}{\partial \boldsymbol{x}} \right\|_F^2. \qquad (14.18)$$
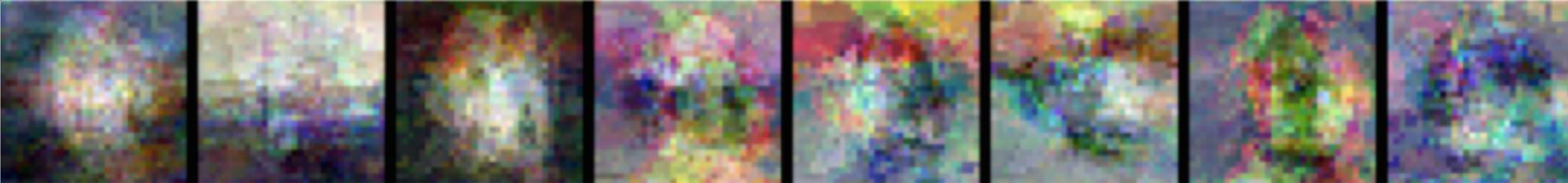


| Input point | Tangent vectors |
|---|---|

Local PCA (no sharing across regions)

Contractive autoencoder

Figure 14.10