

Machine Learning (part II)

Autoencoders

Angelo Ciaramella

Introduction

■ Autoencoder

- neural network that is trained to attempt to copy its input to its output

■ Modern autoencoders

- have generalized the idea of an **encoder** and a **decoder** beyond deterministic functions to stochastic mappings



Principal Component Analysis

- Principal Component Analysis (PCA) is a statistical technique
 - Dimensionality reduction
 - Lossy data compression
 - Feature extraction
 - Data visualization
- It is also known as the *Karhunen-Loeve* transform
- PCA can be defined as the principal subspace such that the variance of the projected data is maximized

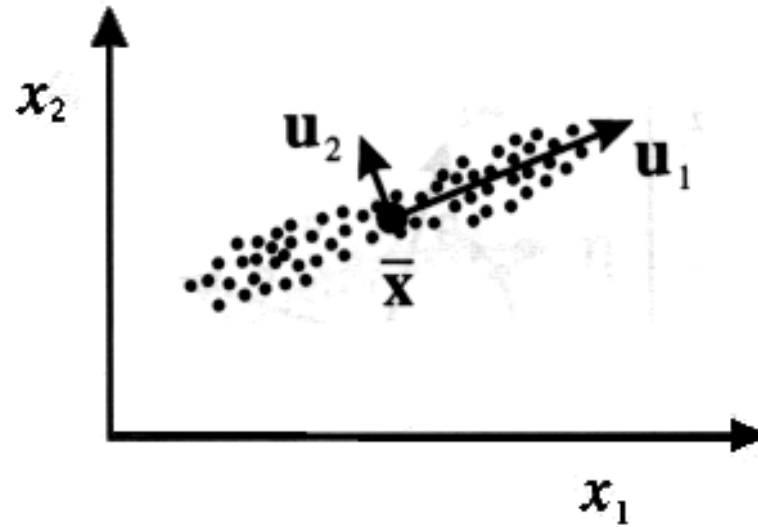


Second-Order methods

- The **second-order methods** are the most popular methods to find a **linear transformation**
- This methods find the representation using only the information contained in the **covariance matrix** of the data vector \mathbf{x}
- **PCA** is widely used in signal processing, statistics, and neural computing



Principal Components



In a linear projection down to one dimension, the optimum choice of projection, in the sense of minimizing the sum-of-squares error, is obtained first subtracting off the mean of the data set, and then projecting onto the first eigenvector \mathbf{u}_1 of the covariance matrix.



Projection error minimization

- We introduce a complete orthonormal set of D -dimensional basis vectors ($i=1, \dots, D$)

$$\mathbf{u}_i^T \mathbf{u}_j = \delta_{ij}$$

- Because this basis is complete, each data point can be represented by a linear combination of the basis vectors

$$\mathbf{x}_n = \sum_{i=1}^D \alpha_{ni} \mathbf{u}_i$$



Projection error minimization

- We can write also that

$$\mathbf{x}_n = \sum_{i=1}^D (\mathbf{x}_n^T \mathbf{u}_i) \mathbf{u}_i \quad \leftarrow \quad \alpha_{nj} = \mathbf{x}_n^T \mathbf{u}_j$$

- Our goal is to **approximate** this data point using a representation involving a restricted number $M < D$ of variables corresponding to a **projection** onto a lower-dimensional subspace

$$\tilde{\mathbf{x}}_n = \sum_{i=1}^M z_{ni} \mathbf{u}_i + \sum_{i=M+1}^D b_i \mathbf{u}_i$$



Projection error minimization

- As our distortion measure we shall use the **squared distance** between the original point and its approximation averaged over the data set so that our goal is to minimize

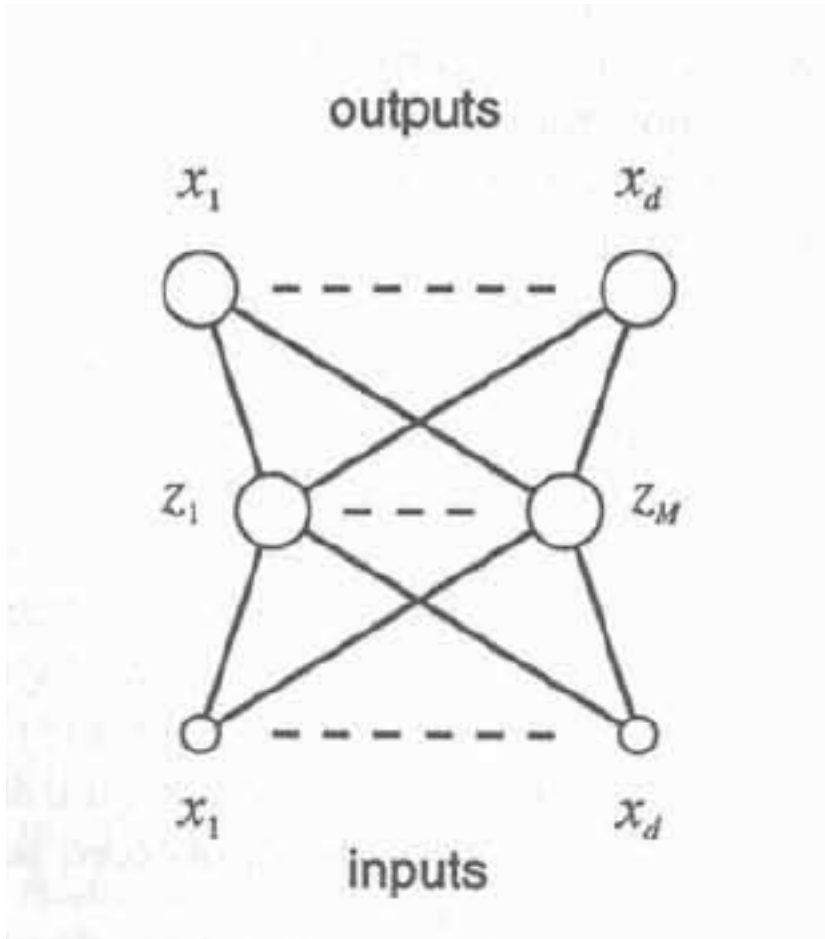
$$J = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \tilde{\mathbf{x}}_n\|^2$$

- The **general solution** is obtained by choosing the basis to be eigenvectors of the covariance matrix given by

$$\mathbf{S}\mathbf{u}_i = \lambda_i \mathbf{u}_i$$



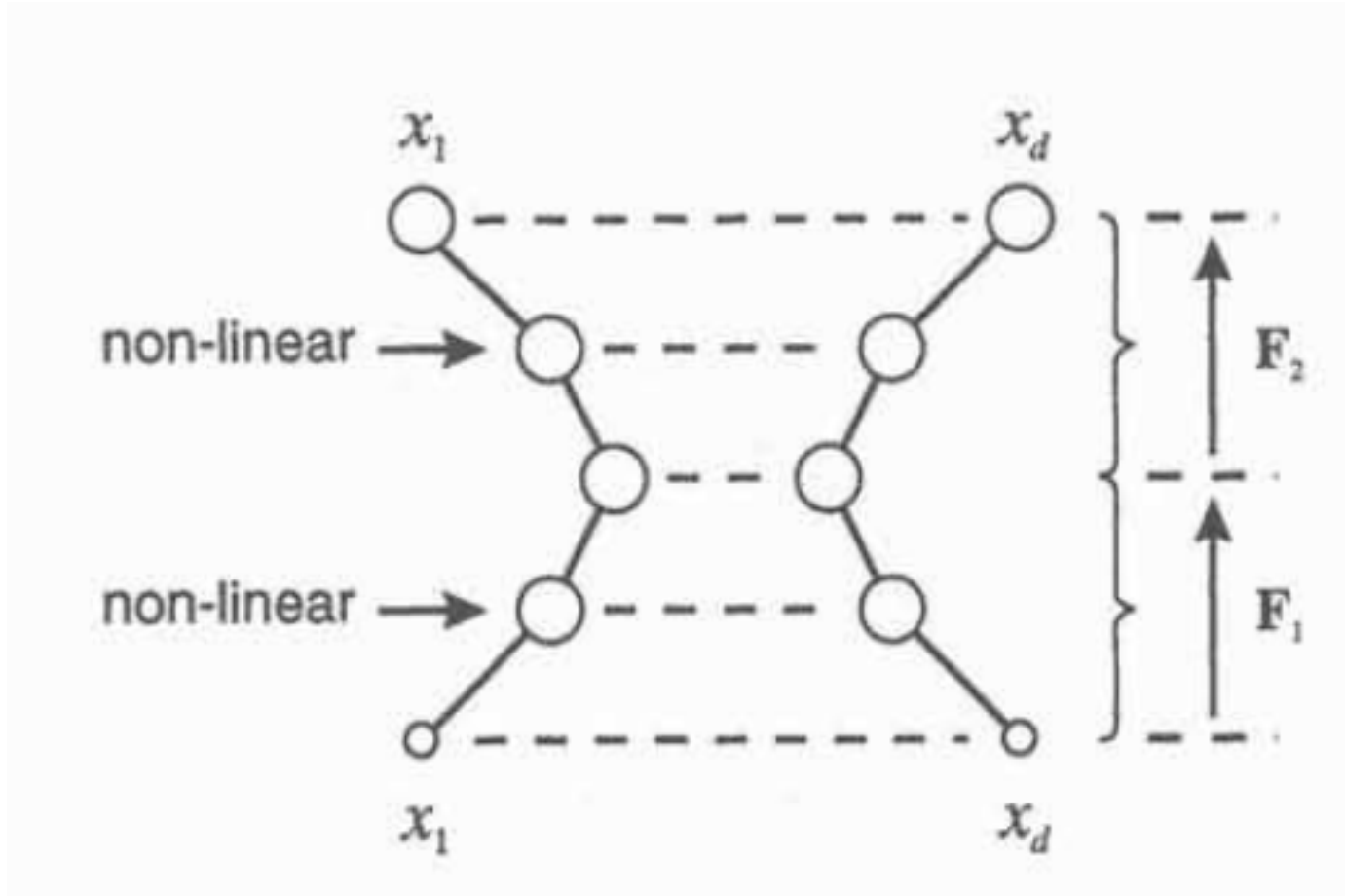
Auto-associative MLP



$$E = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^d \{y_k(\mathbf{x}^n) - x_k^n\}^2.$$



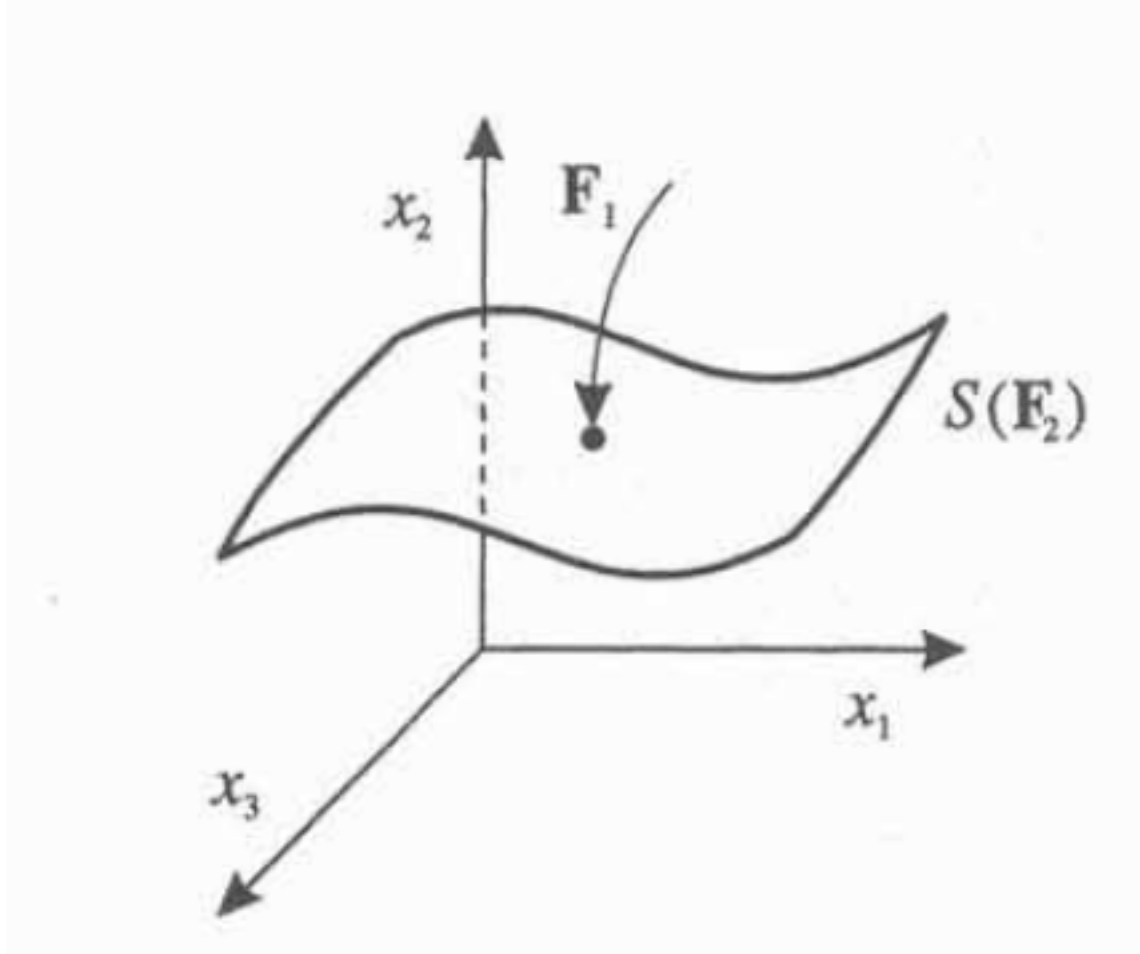
Auto-associative MPL



Non-linear units



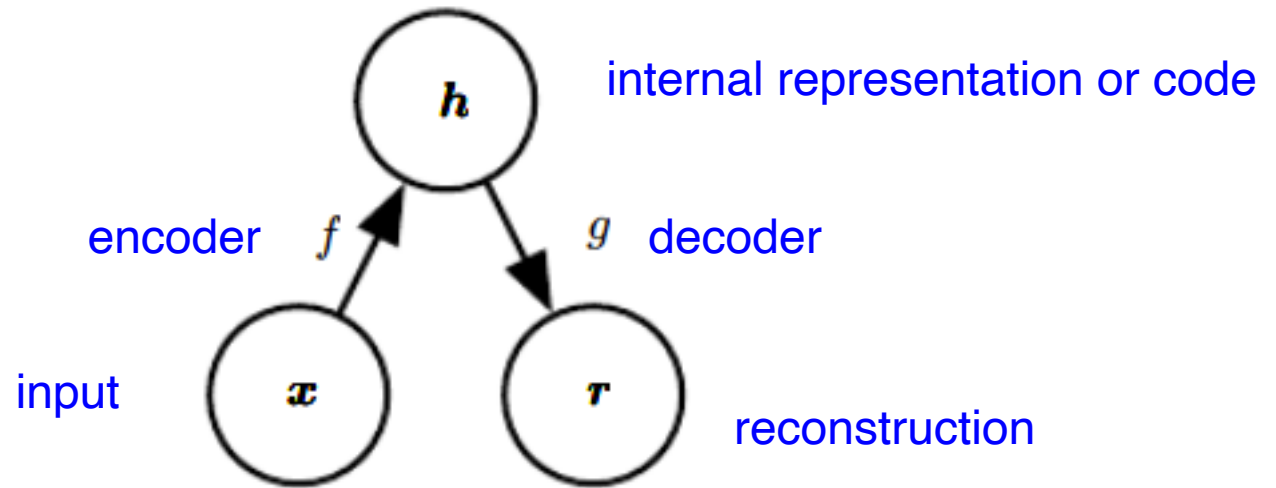
Auto-associative MLP



Mapping of the non-linear auto-associative MLP



Autoencoder



General structure of an autoencoder



Autoencoder

- Undercomplete autoencoder
 - code dimension is less than the input dimension
 - representation forces the autoencoder to capture the most salient features of the training data

- Loss function

$$L(\mathbf{x}, g(f(\mathbf{x})))$$

- PCA subspace
 - decoder linear
 - L is mean squared error



Regularized autoencoder

- Regularized autoencoders
 - use a loss function that encourages the model to have **other properties** besides the ability to copy its input to its output
 - Examples
 - Sparsity of the representation
 - Smallness of the derivative of the representation
 - Robustness to noise
 - Robustness to missing inputs
 - Generative models
 - Helmholtz machine
 - Variational autoencoder



Sparse autoencoder

- Sparse penalty

$$L(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{h})$$

$$\mathbf{h} = f(\mathbf{x}) \quad \text{typically}$$

- used to **learn features** for another task such as classification



Sparse autoencoder

- Joint distribution

$$p_{\text{model}}(\mathbf{x}, \mathbf{h}) = p_{\text{model}}(\mathbf{h}) p_{\text{model}}(\mathbf{x} \mid \mathbf{h})$$

- Log-likelihood

$$\log p_{\text{model}}(\mathbf{x}) = \log \sum_{\mathbf{h}} p_{\text{model}}(\mathbf{h}, \mathbf{x})$$

- Maximizing

$$\log p_{\text{model}}(\mathbf{h}, \mathbf{x}) = \log p_{\text{model}}(\mathbf{h}) + \log p_{\text{model}}(\mathbf{x} \mid \mathbf{h})$$



Sparse autoencoder

- e.g., Laplace prior

$$p_{\text{model}}(h_i) = \frac{\lambda}{2} e^{-\lambda|h_i|}$$

$$\Omega(\mathbf{h}) = \lambda \sum_i |h_i|$$

$$-\log p_{\text{model}}(\mathbf{h}) = \sum_i \left(\lambda|h_i| - \log \frac{\lambda}{2} \right) = \Omega(\mathbf{h}) + \text{const}$$

Way of approximately training a generative model

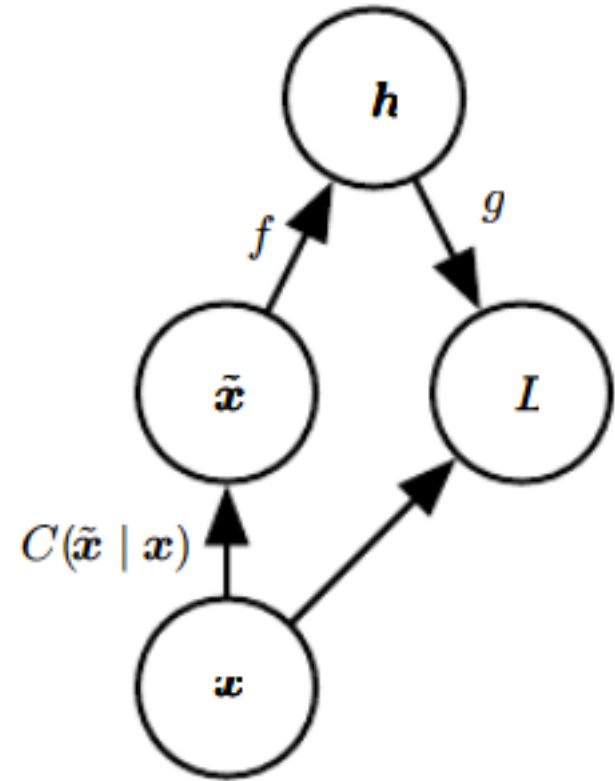


Denoising autoencoder

■ Minimization

$$L(\mathbf{x}, g(f(\tilde{\mathbf{x}})))$$

Noisy version of \mathbf{x}



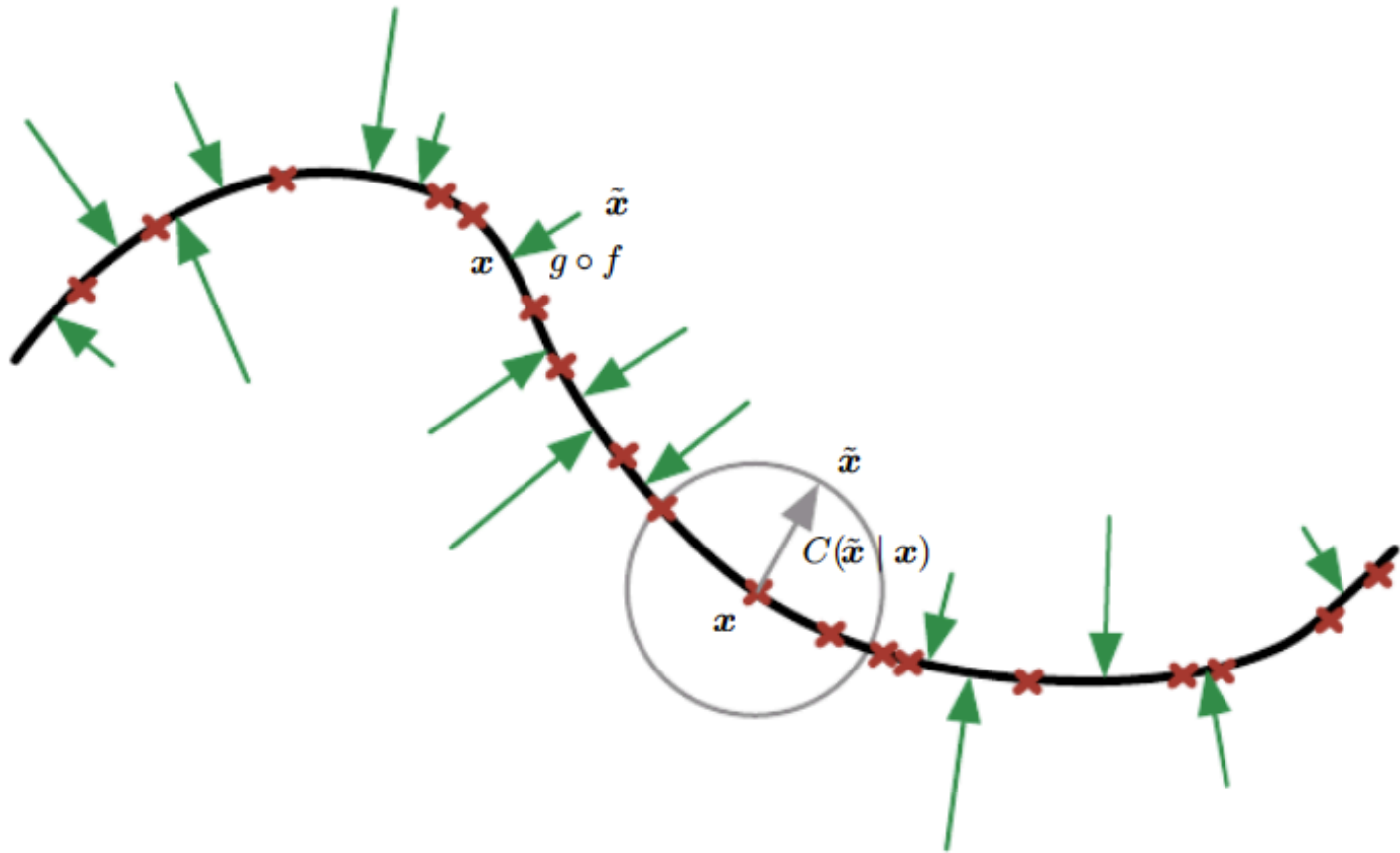
Denoising autoencoder

1. Sample a training example \mathbf{x} from the training data.
2. Sample a corrupted version $\tilde{\mathbf{x}}$ from $C(\tilde{\mathbf{x}} | \mathbf{x} = \mathbf{x})$.
3. Use $(\mathbf{x}, \tilde{\mathbf{x}})$ as a training example for estimating the autoencoder reconstruction distribution $p_{\text{reconstruct}}(\mathbf{x} | \tilde{\mathbf{x}}) = p_{\text{decoder}}(\mathbf{x} | \mathbf{h})$ with \mathbf{h} the output of encoder $f(\tilde{\mathbf{x}})$ and p_{decoder} typically defined by a decoder $g(\mathbf{h})$.

$$- \mathbb{E}_{\mathbf{x} \sim \hat{p}_{\text{data}}(\mathbf{x})} \mathbb{E}_{\tilde{\mathbf{x}} \sim C(\tilde{\mathbf{x}} | \mathbf{x})} \log p_{\text{decoder}}(\mathbf{x} | \mathbf{h} = f(\tilde{\mathbf{x}}))$$



Denoising autoencoder



denoising autoencoder is trained to map a corrupted data



Penalizing Derivatives

■ Minimization

contractive autoencoder

$$L(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{h}, \mathbf{x})$$

$$\Omega(\mathbf{h}, \mathbf{x}) = \lambda \sum_i \|\nabla_{\mathbf{x}} h_i\|^2$$

Forces the autoencoder to learn features that capture information about the training distribution



Depth

- Depth
 - using deep encoders and decoders offers many advantages
 - One major advantage of non-trivial depth is that the universal approximator theorem guarantees
 - yield much better compression than corresponding shallow or linear autoencoders

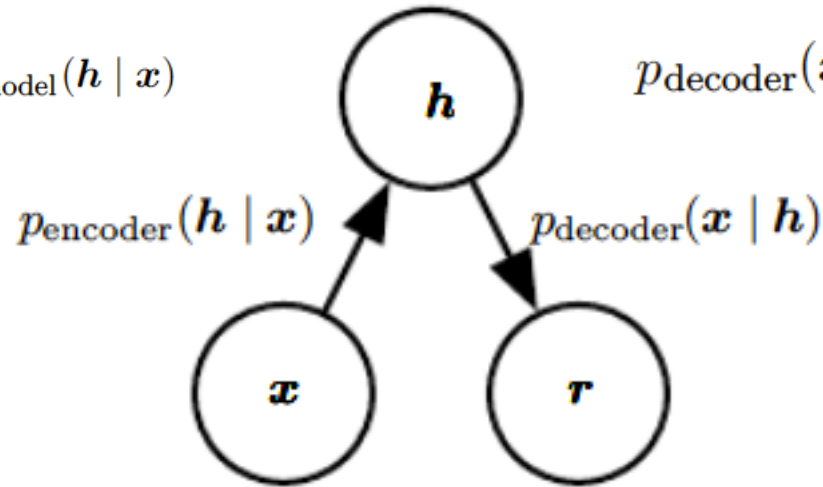


Stochastic autoencoder

■ Depth

$$p_{\text{encoder}}(\mathbf{h} \mid \mathbf{x}) = p_{\text{model}}(\mathbf{h} \mid \mathbf{x})$$

$$p_{\text{decoder}}(\mathbf{x} \mid \mathbf{h}) = p_{\text{model}}(\mathbf{x} \mid \mathbf{h})$$



Both the encoder and the decoder are not simple functions but instead involve some **noise injection**, meaning that their output can be seen as sampled from a distribution



Contractive autoencoders

- Encouraging the derivatives of f to be as small as possible

Jacobian matrix

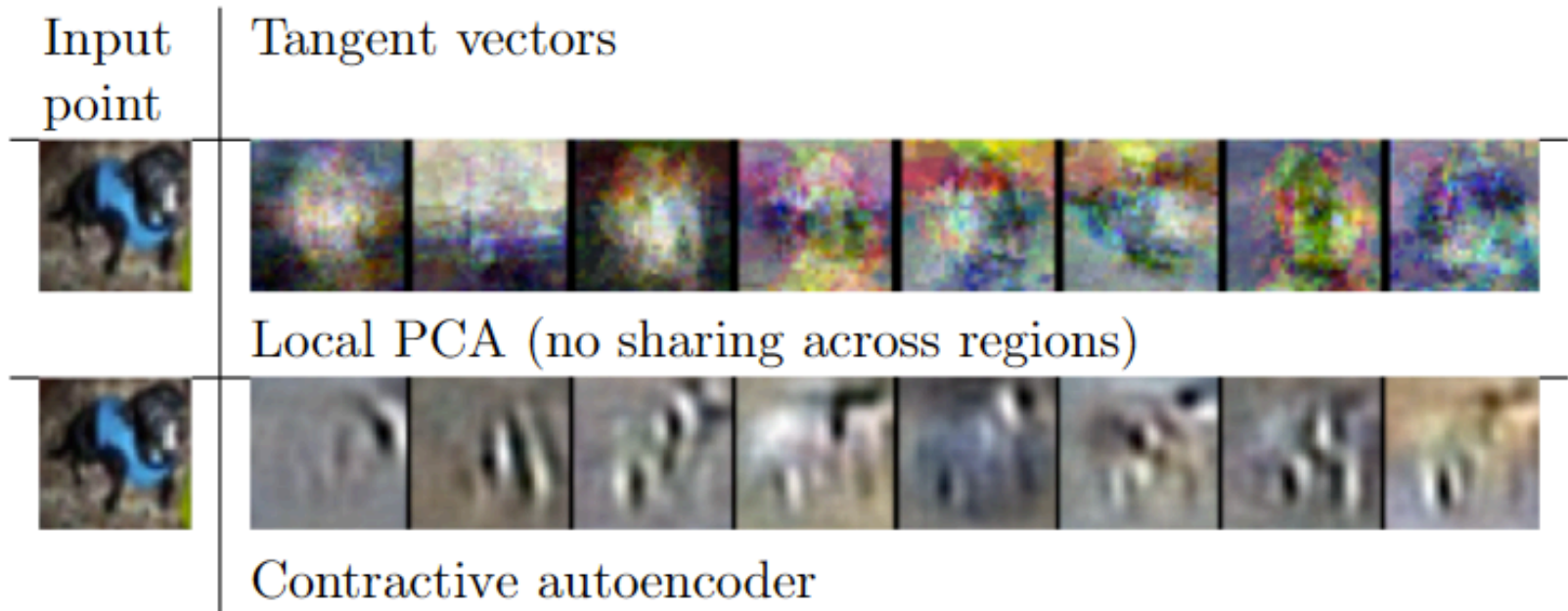
$$\Omega(\mathbf{h}) = \lambda \left\| \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right\|_F^2$$

Frobenius norm (sum of squared elements)



Contractive autoencoders

largest singular values are interpreted as the tangent directions that the contractive autoencoder has learned



learn tangent vectors that show how the image changes as objects in the image gradually change pose

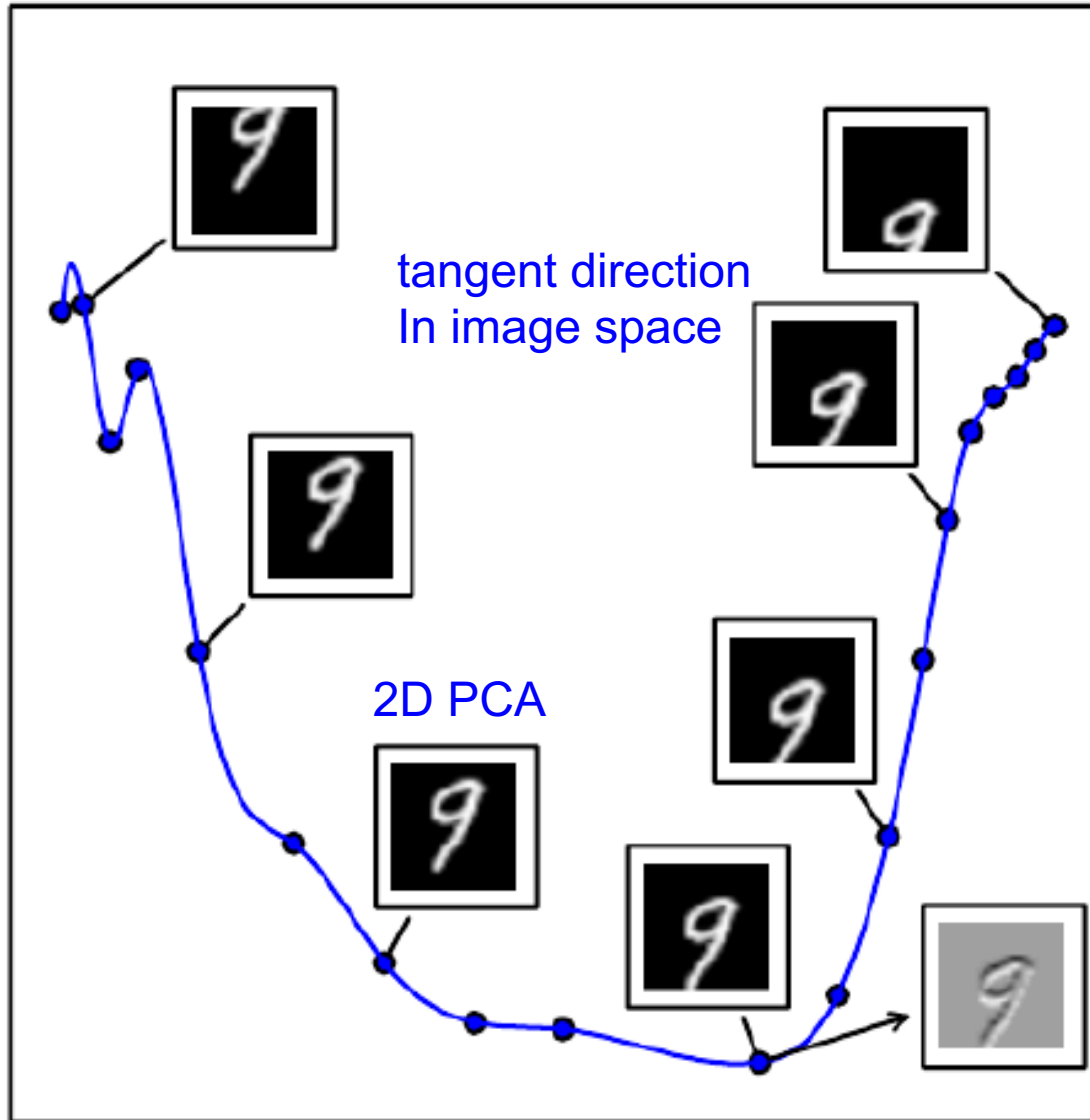


Manifold learning

- data concentrates around a low-dimensional manifold
- tangent planes
 - data concentrates around a low-dimensional manifold
 - At a point x on a d -dimensional manifold, the tangent plane is given by d **basis vectors** that span the local directions of **variation allowed on the manifold**



Tangent plane



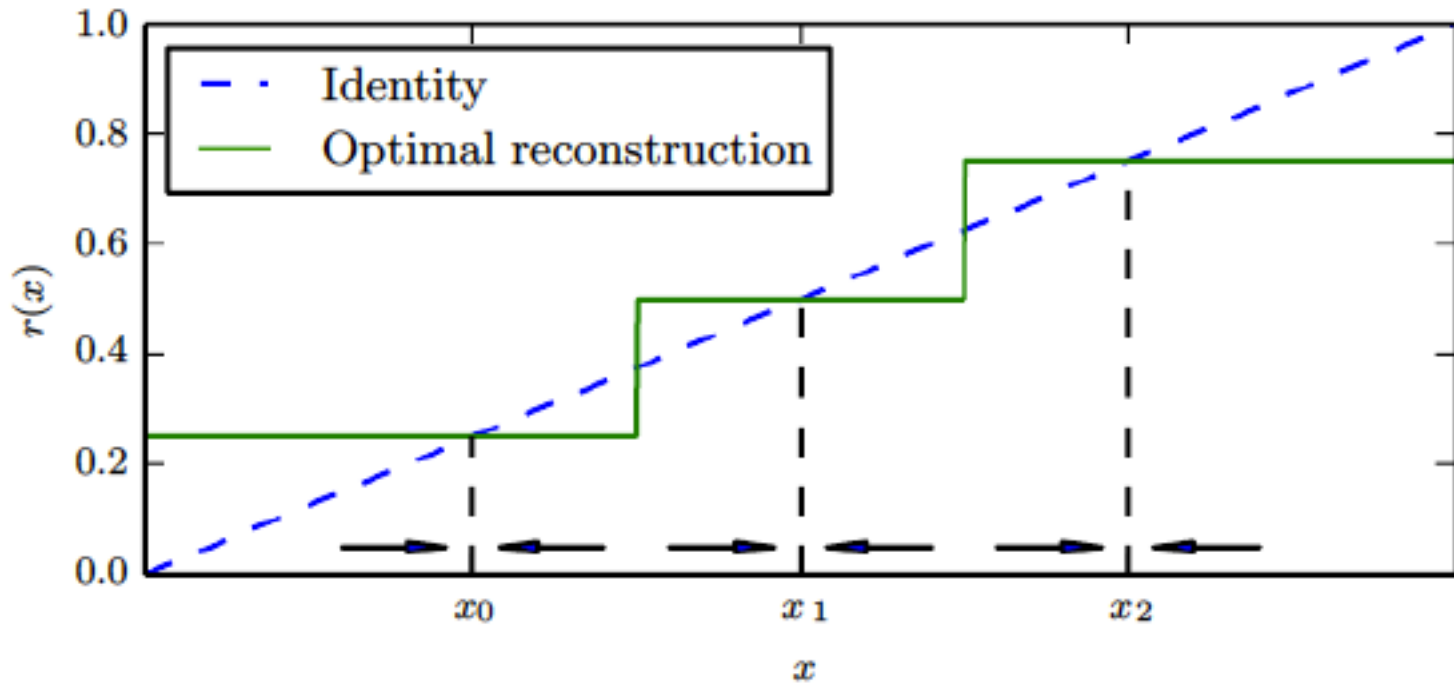
one-dimensional manifold
in 784-dimensional space

translating it vertically

vertical translation
defines a coordinate
along a one-dimensional
manifold



Reconstruction



If the autoencoder learns a reconstruction function that is invariant to small perturbations near the data points, it captures the manifold structure of the data.



Variational Autoencoders

■ Inference as optimization

- Approximate inference algorithms may then be derived by approximating the underlying optimization problem

- Aim

$$\log p(\mathbf{v}; \boldsymbol{\theta})$$

■ Evidence lower bound (ELBO)

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \log p(\mathbf{v}; \boldsymbol{\theta}) - D_{\text{KL}}(q(\mathbf{h} | \mathbf{v}) || p(\mathbf{h} | \mathbf{v}; \boldsymbol{\theta}))$$

\mathbf{v}
observed variables

\mathbf{h}
latent variables

q
arbitrary probability distribution



Variational Autoencoders

■ Canonical definition

$$\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) = \mathbb{E}_{\mathbf{h} \sim q} [\log p(\mathbf{h}, \mathbf{v})] + H(q)$$

$$\begin{aligned}\mathcal{L}(\mathbf{v}, \boldsymbol{\theta}, q) &= \log p(\mathbf{v}; \boldsymbol{\theta}) - D_{\text{KL}}(q(\mathbf{h} | \mathbf{v}) \| p(\mathbf{h} | \mathbf{v}; \boldsymbol{\theta})) \\ &= \log p(\mathbf{v}; \boldsymbol{\theta}) - \mathbb{E}_{\mathbf{h} \sim q} \log \frac{q(\mathbf{h} | \mathbf{v})}{p(\mathbf{h} | \mathbf{v})} \\ &= \log p(\mathbf{v}; \boldsymbol{\theta}) - \mathbb{E}_{\mathbf{h} \sim q} \log \frac{q(\mathbf{h} | \mathbf{v})}{\frac{p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta})}{p(\mathbf{v}; \boldsymbol{\theta})}} \\ &= \log p(\mathbf{v}; \boldsymbol{\theta}) - \mathbb{E}_{\mathbf{h} \sim q} [\log q(\mathbf{h} | \mathbf{v}) - \log p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta}) + \log p(\mathbf{v}; \boldsymbol{\theta})] \\ &= - \mathbb{E}_{\mathbf{h} \sim q} [\log q(\mathbf{h} | \mathbf{v}) - \log p(\mathbf{h}, \mathbf{v}; \boldsymbol{\theta})].\end{aligned}$$



Variational Autoencoders

■ Variational Autoencoders

$p_{\text{model}}(\mathbf{z})$

distribution

$g(\mathbf{z})$

generator network

joint log-likelihood

$$\begin{aligned}\mathcal{L}(q) &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{z}, \mathbf{x}) + \mathcal{H}(q(\mathbf{z} | \mathbf{x})) \\ &= \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z}|\mathbf{x})} \log p_{\text{model}}(\mathbf{x} | \mathbf{z}) - D_{\text{KL}}(q(\mathbf{z} | \mathbf{x}) || p_{\text{model}}(\mathbf{z})) \\ &\leq \log p_{\text{model}}(\mathbf{x}).\end{aligned}$$



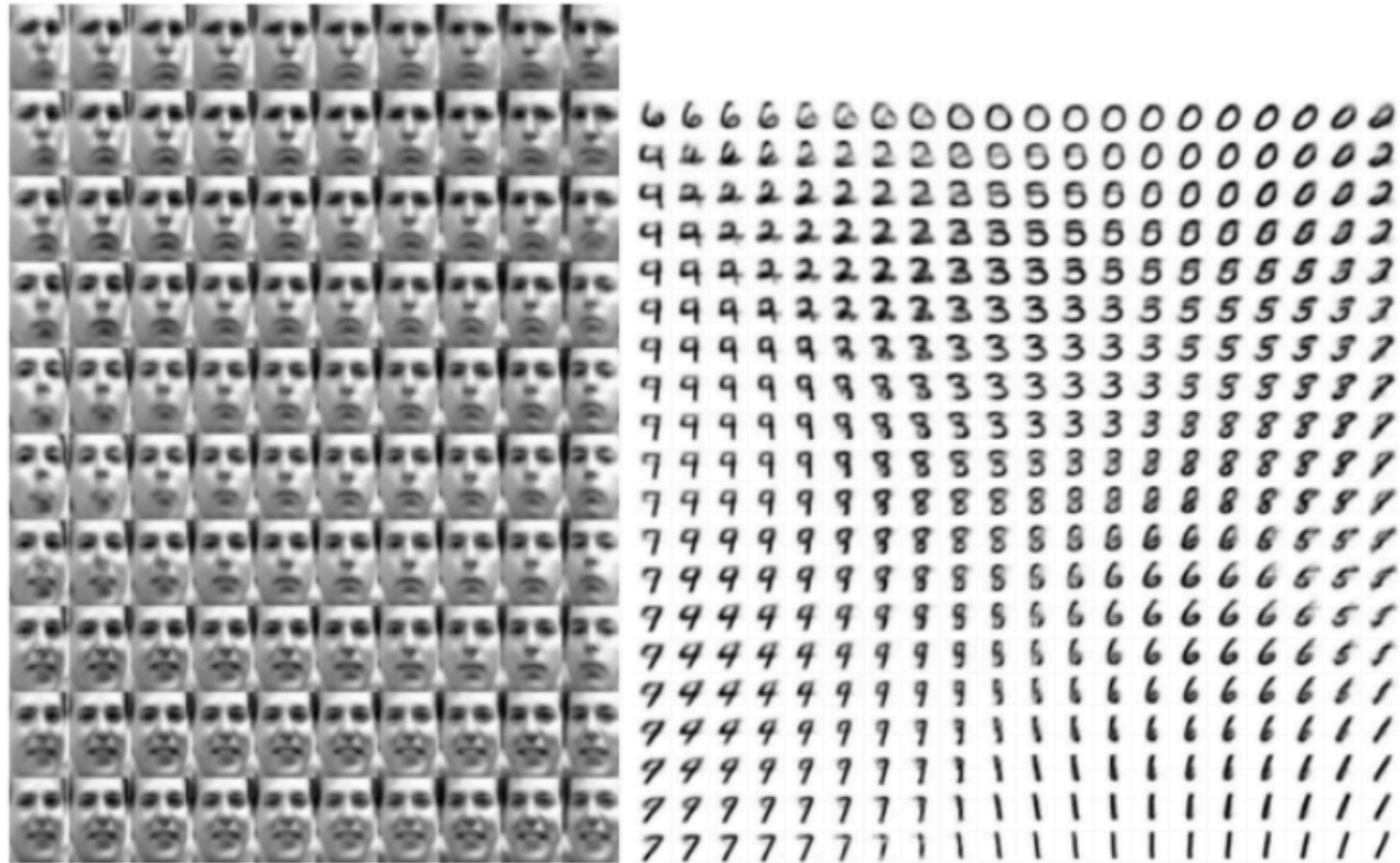
Variational Autoencoders

■ Training

- train a **parametric encoder** that produces the parameters of q
- maximizing L with respect to the parameters of the encoder and decoder
- All of the expectations in L may be approximated by **Monte Carlo sampling**



Interesting properties



algorithm discovered two independent factors of variation present in images of faces: angle of rotation and emotional expression

