

Chapter 14

Autoencoders

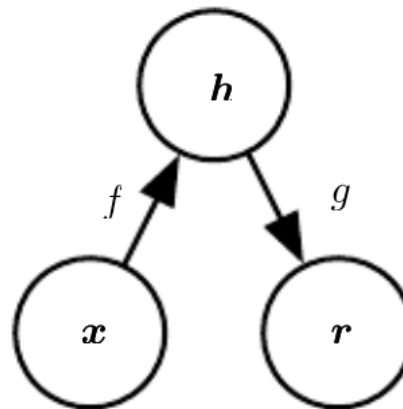
Autoencoders

- A type of neural networks trained to copy **approximately** its input to its output in the hopes of learning useful features
- The network of an autoencoder may be viewed as containing an encoder and a decoder, specifying deterministic or stochastic mappings

Encoder: $\mathbf{h} = f(\mathbf{x})$ or $p_{\text{model}}(\mathbf{h}|\mathbf{x})$

Decoder: $\mathbf{r} = g(\mathbf{h})$ or $p_{\text{model}}(\mathbf{x}|\mathbf{h})$

where the hidden layer \mathbf{h} describes a code used to represent \mathbf{x}



- The learning is to minimize a loss function, likely with regularization

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

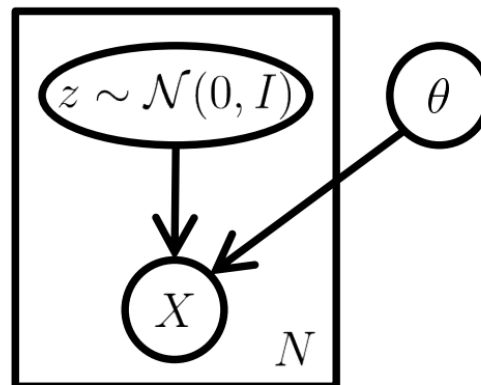
- Most learning techniques for training feedforward networks can apply
- Traditionally, autoencoders were used for dimension reduction
- However, theoretical connections between autoencoders and some modern latent variable models have brought autoencoders to the forefront of generative modeling

Variational Autoencoders (VAE)

- A probabilistic generative model with latent variables that is built on top of end-to-end trainable neural networks

$$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$$

$$p(\mathbf{x}|\mathbf{z}) = \underbrace{p(\mathbf{x}; o(\mathbf{z}; \boldsymbol{\theta}))}_{\text{Neural Networks}} = \mathcal{N}(\mathbf{x}; o(\mathbf{z}; \boldsymbol{\theta}), \sigma^2 \mathbf{I})$$



Training VAE

- To determine θ , we would intuitively hope to maximize the marginal distribution $p(\mathbf{x}; \theta)$

$$p(\mathbf{x}; \theta) = \int p(\mathbf{x}|\mathbf{z}; \theta)p(\mathbf{z})d\mathbf{z}$$

- This however becomes difficult as the integration over \mathbf{z} is in general intractable when $p(\mathbf{x}|\mathbf{z}; \theta)$ is modeled by a neural network
- To circumvent this difficulty, we recall that

$$\log p(\mathbf{X}; \theta) = \mathcal{L}(\mathbf{X}, q, \theta) + \text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}; \theta))$$

where

$$\mathcal{L}(\mathbf{X}, q, \theta) = \int q(\mathbf{Z}) \log p(\mathbf{X}, \mathbf{Z}; \theta) d\mathbf{Z} - \int q(\mathbf{Z}) \log q(\mathbf{Z}) d\mathbf{Z}$$

$$\text{KL}(q(\mathbf{Z})||p(\mathbf{Z}|\mathbf{X}; \theta)) = \int q(\mathbf{Z}) \log \frac{q(\mathbf{Z})}{p(\mathbf{Z}|\mathbf{X}; \theta)} d\mathbf{Z}$$

- A rearrangement gives

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

- As the equality holds for any choice of $q(\mathbf{Z})$, we introduce a distribution $q(\mathbf{Z} | \mathbf{X}; \boldsymbol{\theta}')$ modeled by another neural network with parameter $\boldsymbol{\theta}'$ to obtain

$$\log p(\mathbf{X}; \boldsymbol{\theta}) - \text{KL}(q(\mathbf{Z} | \mathbf{X}; \boldsymbol{\theta}') || p(\mathbf{Z} | \mathbf{X}; \boldsymbol{\theta})) = \mathcal{L}(\mathbf{X}, q, \boldsymbol{\theta})$$

- The right hand side can be spell out as

$$\begin{aligned} \mathcal{L}(\mathbf{X}, q, \boldsymbol{\theta}) &= E_{\mathbf{Z} \sim q(\mathbf{Z} | \mathbf{X}; \boldsymbol{\theta}')} p(\mathbf{X} | \mathbf{Z}; \boldsymbol{\theta}) + E_{\mathbf{Z} \sim q(\mathbf{Z} | \mathbf{X}; \boldsymbol{\theta}')} p(\mathbf{Z}) \\ &\quad - E_{\mathbf{Z} \sim q(\mathbf{Z} | \mathbf{X}; \boldsymbol{\theta}')} q(\mathbf{Z} | \mathbf{X}; \boldsymbol{\theta}') \\ &= E_{\mathbf{Z} \sim q(\mathbf{Z} | \mathbf{X}; \boldsymbol{\theta}')} p(\mathbf{X} | \mathbf{Z}; \boldsymbol{\theta}) - \text{KL}(q(\mathbf{Z} | \mathbf{X}; \boldsymbol{\theta}') || p(\mathbf{Z})) \end{aligned}$$

- Now, instead of directly maximizing the intractable $p(\mathbf{X}; \boldsymbol{\theta})$, we attempt to maximize

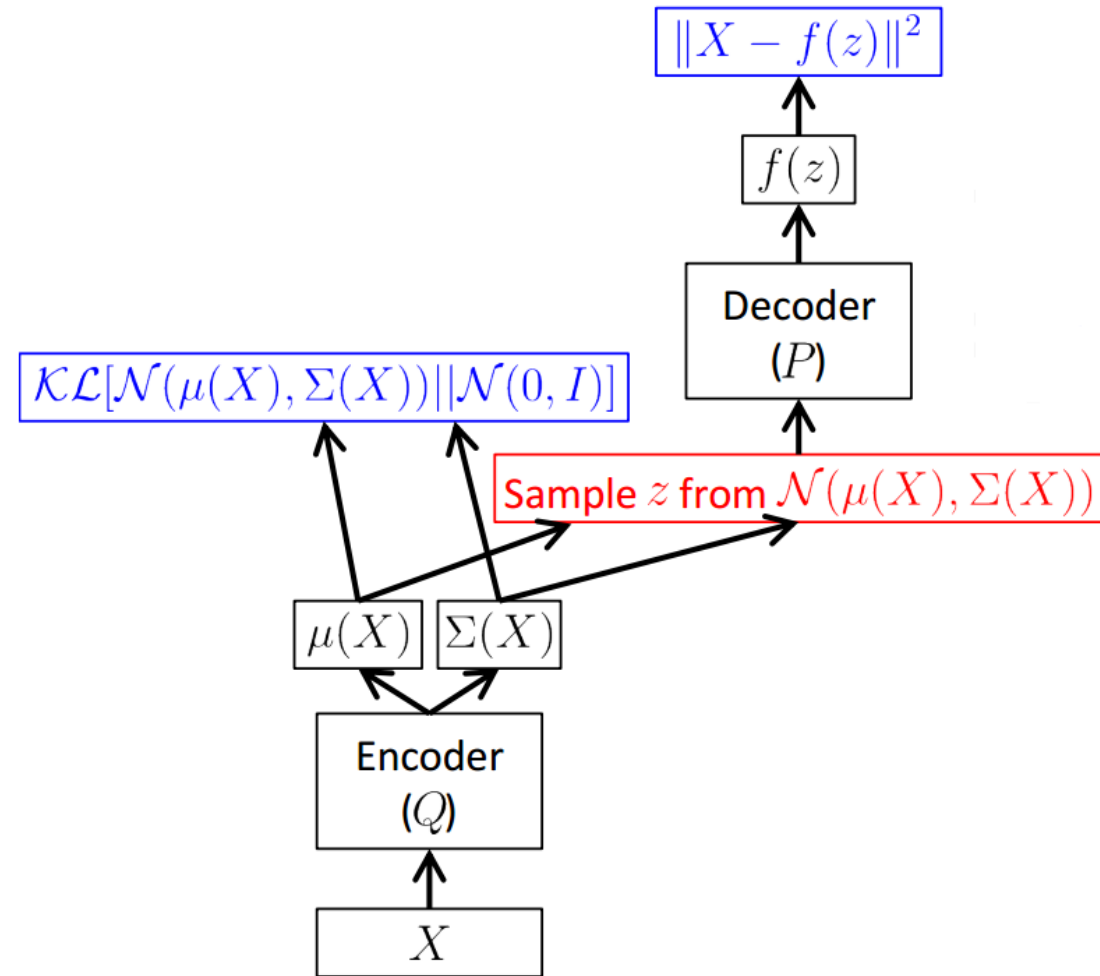
$$\log p(\mathbf{X}; \boldsymbol{\theta}) - \text{KL}(q(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}') || p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}))$$

which amounts to maximizing

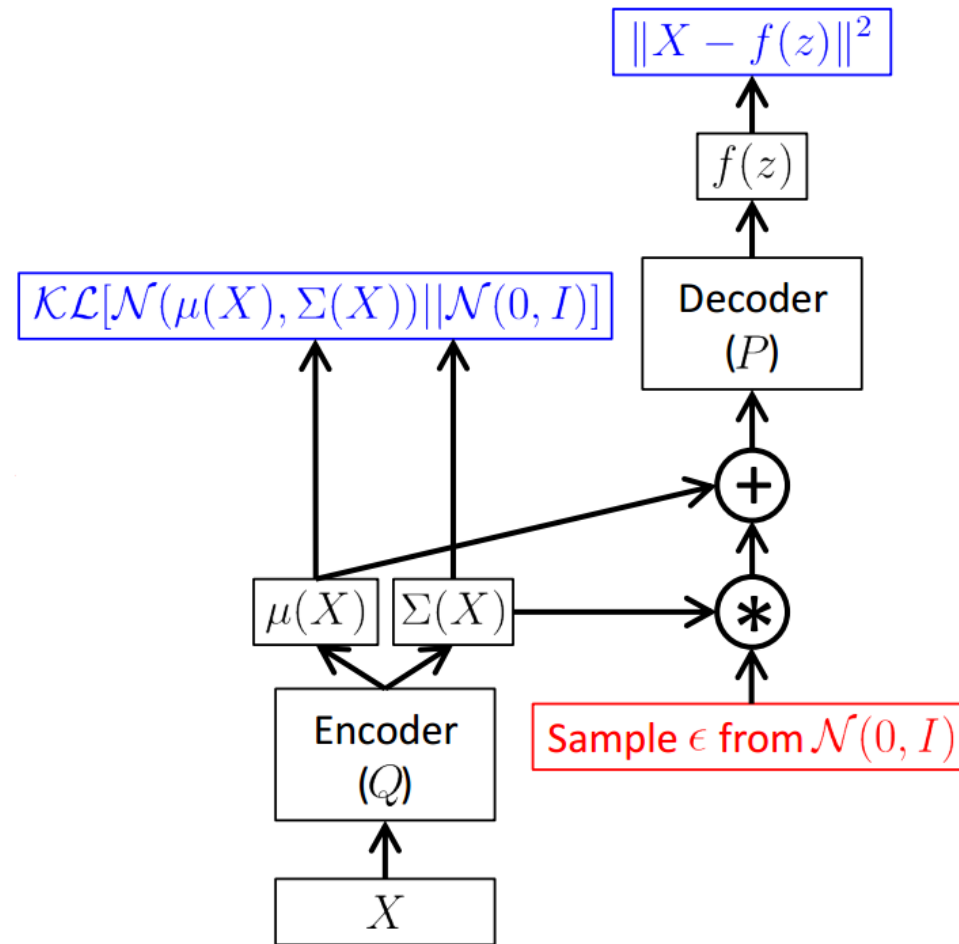
$$E_{\mathbf{Z} \sim q(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}')} p(\mathbf{X}|\mathbf{Z}; \boldsymbol{\theta}) - \text{KL}(q(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}') || p(\mathbf{Z}))$$

- A by-product of this training process is a stochastic encoder

$$q(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}') \approx p(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta})$$



$$\underbrace{E_{\mathbf{Z} \sim q(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}')} p(\mathbf{X}|\mathbf{Z}; \boldsymbol{\theta})}_{\text{Sampling needed}} - \text{KL}(q(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}') || p(\mathbf{Z}))$$



$$\underbrace{E_{\mathbf{Z} \sim q(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}')} p(\mathbf{X}|\mathbf{Z}; \boldsymbol{\theta})}_{\text{Re-parameterization for end-to-end training}} - \text{KL}(q(\mathbf{Z}|\mathbf{X}; \boldsymbol{\theta}') || p(\mathbf{Z}))$$

- Even though the following term can be evaluated by sampling \mathbf{Z} from $q(\mathbf{Z}|\mathbf{X};\boldsymbol{\theta}')$, it becomes difficult to compute the gradient w.r.t. $\boldsymbol{\theta}'$

$$E_{\mathbf{Z} \sim q(\mathbf{Z}|\mathbf{X};\boldsymbol{\theta}')} p(\mathbf{X}|\mathbf{Z};\boldsymbol{\theta})$$

- The re-parameterization technique works around this difficulty by generating samples input to the decoder with

$$\mathbf{B}(\mathbf{X})\boldsymbol{\epsilon} + \boldsymbol{\mu}(\mathbf{X})$$

where $\mathbf{B}\mathbf{B}^T = \Sigma$ and $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$

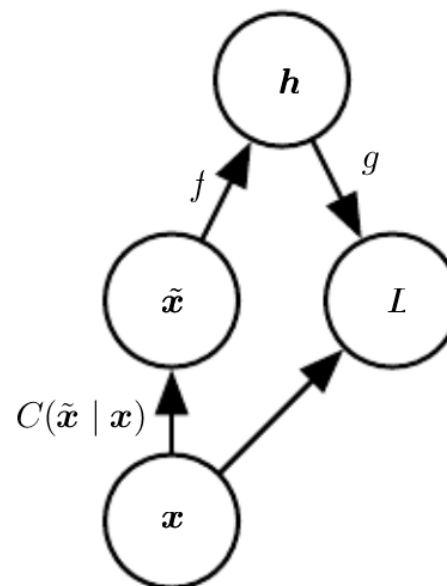
- In fact, the encoder can learn $\mathbf{B}(\mathbf{X})$ directly

Denoising Autoencoders (DAE)

- The DAE is to receive a corrupted data point as input and to predict the uncorrupted data point as output; that is, to minimize

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

where $\tilde{\mathbf{x}}$ is a noise-corrupted version of \mathbf{x}



- To be precise, the training of DAE proceeds as follows
 1. Sample an x from the training data
 2. Sample a corrupted version \tilde{x} from $C(\tilde{x}|x)$
 3. Minimize the negative log-likelihood by performing gradient descent w.r.t. model parameters

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

- When the encoder f is deterministic, the training of DAE is no different than training a feedforward network

- It is shown that when both $p_{\text{decoder}}(\mathbf{x}|\mathbf{h})$ and $C(\tilde{\mathbf{x}}|\mathbf{x})$ are assumed to be Gaussian, i.e., training with

$$\min \|g(f(\tilde{\mathbf{x}})) - \mathbf{x}\|^2 \text{ and } C(\tilde{\mathbf{x}}|\mathbf{x}) \sim \mathcal{N}(\tilde{\mathbf{x}}; \mathbf{x}, \sigma^2 \mathbf{I}),$$

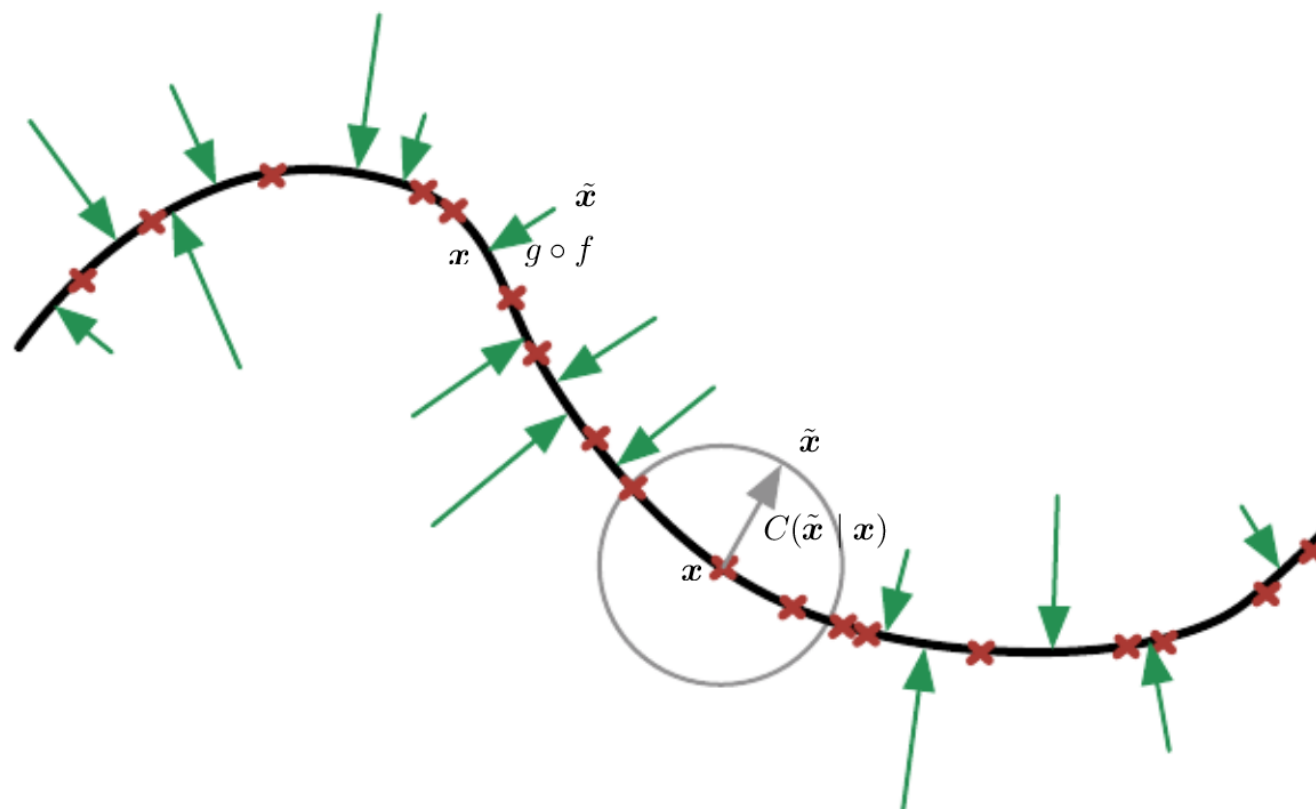
the DAE learns a vector field $(g(f(\mathbf{x})) - \mathbf{x})$ that gives estimates of **the score of the data distribution** defined as

$$\nabla_{\mathbf{x}} \log p(\mathbf{x})$$

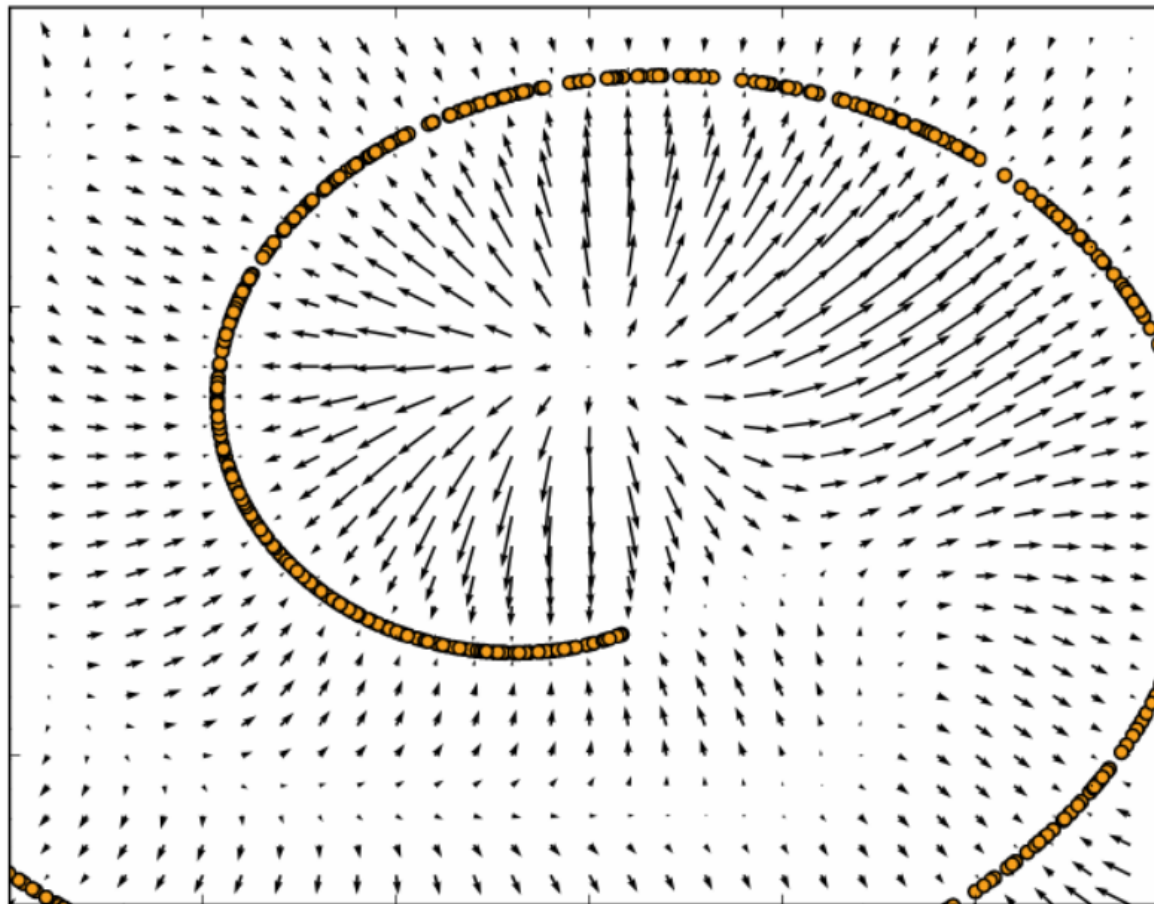
- Note that when $\|g(f(\tilde{\mathbf{x}})) - \mathbf{x}\|^2$ is minimized, we have

$$g(f(\tilde{\mathbf{x}})) \approx E_{\mathbf{x}, \tilde{\mathbf{x}} \sim \hat{p}_{\text{data}}(\mathbf{x})C(\tilde{\mathbf{x}}|\mathbf{x})}[\mathbf{x}|\tilde{\mathbf{x}}]$$

- Thus, $(g(f(\tilde{\mathbf{x}})) - \tilde{\mathbf{x}})$ is a vector that points approximately back to the nearest point on the data manifold



Green arrows: $g(f(\tilde{x})) - \tilde{x}$



Vector field learned by a DAE
(Vector field has zeros at both maxima and minima of $p(\mathbf{x})$)

Sparse Autoencoders

- A sparse autoencoder is an autoencoder whose training criterion involves a sparsity penalty $\Omega(\mathbf{h})$

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

- It can be interpreted as approximating the maximum likelihood training of a generative model $p_{\text{model}}(\mathbf{x}, \mathbf{h})$ with latent variables \mathbf{h}

$$\begin{aligned} \log p_{\text{model}}(\mathbf{x}) &= \log \sum_{\mathbf{h}} p_{\text{model}}(\mathbf{x}, \mathbf{h}) \\ &\approx \underbrace{\log p_{\text{model}}(\mathbf{h})}_{\Omega} + \underbrace{\log p_{\text{model}}(\mathbf{x}|\mathbf{h})}_{L}, \end{aligned}$$

where the $p_{\text{model}}(\mathbf{h})$ is factorial and follows the Laplace prior

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

Contractive Autoencoders (CAE)

- The CAE imposes a regularizer on the code \mathbf{h} which encourages to learn an encoder function that does not change much when input \mathbf{x} changes slightly

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

where

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

- The encoder $f(\mathbf{x})$ at a training point \mathbf{x}_0 can be approximated as

$$f(\mathbf{x}) \approx f(\mathbf{x}_0) + \frac{\partial f(\mathbf{x}_0)}{\partial \mathbf{x}} (\mathbf{x} - \mathbf{x}_0)$$

- As such, the CAE is seen to encourage the Jacobian matrix $\partial f(\mathbf{x}_0)/\partial \mathbf{x}$ at every training point \mathbf{x}_0 to become contractive, making their singular values become as small as possible

- It however is noticed that the optimization has to respect also the reconstruction error; this leads to an effect that keeps the singular values along directions with large local variances
- These directions are known as **tangent directions** to the data manifold; that is, they correspond to real variations in the data

Review

- Stochastic vs. deterministic autoencoders
- Autoencoders vs. generative models with latent variables
- Training autoencoders vs. learning data manifolds