

## Autoencoders

### A: The autoencoder formalism

## Unnamed structure

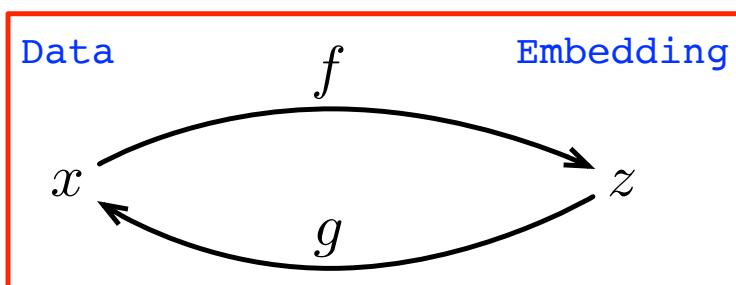
So far: representations that capture various types of structure:

- clusters
- linear subspaces
- manifolds
- explicit distance and similarity information

Can we find representations that capture:

- **multiple types of structure**, simultaneously
- **other structure for which we don't even have definitions?**

## Deterministic autoencoders



- $f$ : encoding map
- $g$ : decoding map

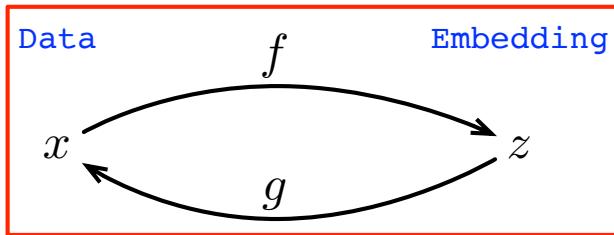
$z$  should be *simpler* than  $x$ :

- e.g., lower-dimensional or sparser
- expose key structure in  $x$

Typically: find  $f, g$  to minimize reconstruction loss  $\mathbb{E}_X[\ell(X, g(f(X)))]$

## PCA as an autoencoder

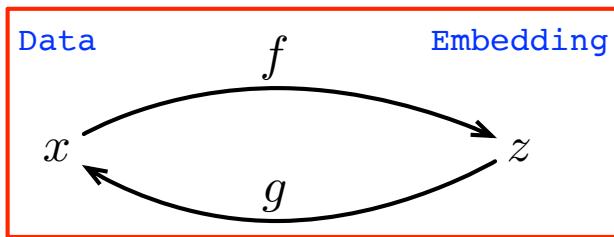
Use PCA to map data from  $\mathbb{R}^d$  to  $\mathbb{R}^k$ .



- What are  $f$  and  $g$ ?
- What loss function is being minimized?

## $K$ -means as an autoencoder

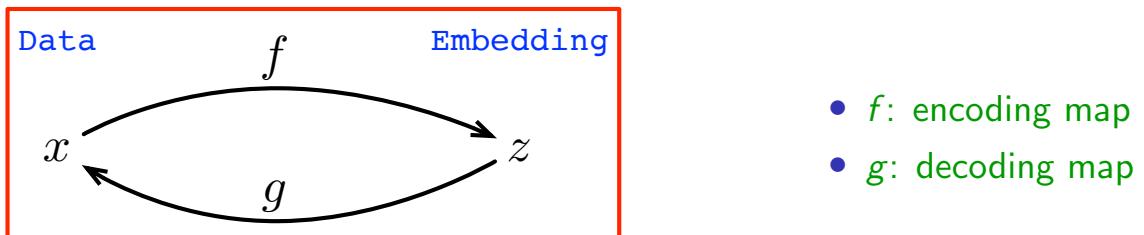
Apply  $k$ -means clustering to data in  $\mathbb{R}^d$ .



- What are  $f$  and  $g$ ?
- What loss function is being minimized?

## More general autoencoders

Allow  $f, g$  to be arbitrary functions, e.g. neural net with one hidden layer.

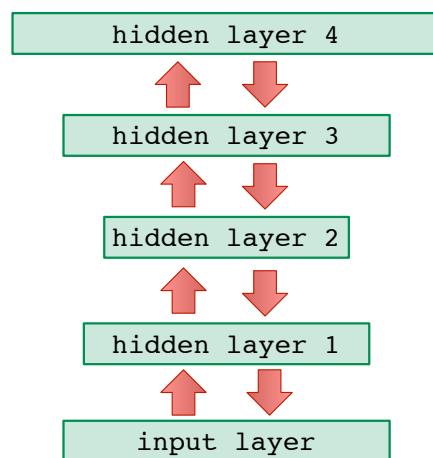


Common variant: **sparse autoencoder**.  $z$  is sparse, e.g.

- by explicitly keeping just the  $k$  highest activations
- by adding a sparsity-encouraging penalty term in the optimization

## Stacked autoencoders

Successively higher-level representations



One way to fit these models:

- Fit one layer at a time to the previous layer's activations
- Then fine-tune the whole structure to minimize reconstruction error

## B: Special case: Dictionary learning

### Linear decoding: dictionary learning

Given data points  $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$ , and an integer  $m$ :

- Choose  $m$  dictionary vectors  $\phi_1, \dots, \phi_m \in \mathbb{R}^d$ .
- Approximate each  $x^{(i)}$  by a linear combination of these dictionary elements.

$$\underbrace{\begin{pmatrix} x^{(1)} & \dots & x^{(n)} \end{pmatrix}}_{\text{data matrix } X} \approx \underbrace{\begin{pmatrix} \phi_1 & \phi_2 & \dots & \phi_m \end{pmatrix}}_{\text{dictionary } \Phi} \underbrace{\begin{pmatrix} z^{(1)} & \dots & z^{(n)} \end{pmatrix}}_{\text{encoding } Z}$$

- **Principal component analysis:** the  $\phi_i$  are orthogonal and  $m \leq d$
- **Nonnegative matrix factorization:**  $\Phi$  and  $Z$  are non-negative
- **K-means:**  $m = k$  and the encodings are in  $\{e_1, \dots, e_k\}$

## More dictionary learning

Given data points  $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$ , and an integer  $m$ :

- Choose  $m$  dictionary vectors  $\phi_1, \dots, \phi_m \in \mathbb{R}^d$ .
- Approximate each  $x^{(i)}$  by a linear combination of these dictionary elements.

$$\underbrace{\begin{pmatrix} x^{(1)} & \dots & x^{(n)} \end{pmatrix}}_{\text{data matrix } X} \approx \underbrace{\begin{pmatrix} \phi_1 & \phi_2 & \dots & \phi_m \end{pmatrix}}_{\text{dictionary } \Phi} \underbrace{\begin{pmatrix} z^{(1)} & \dots & z^{(n)} \end{pmatrix}}_{\text{encoding } Z}$$

- **Independent component analysis:** the rows of  $Z$  are (approximately) statistically independent and  $m \leq d$
- **Sparse coding:** columns of  $Z$  are sparse and often  $m > d$  ("overcomplete basis")

## Sparse coding

Given  $x^{(1)}, \dots, x^{(n)} \in \mathbb{R}^d$ , find dictionary vectors  $\phi_1, \dots, \phi_m$  and sparse representations  $z^{(1)}, \dots, z^{(n)} \in \mathbb{R}^m$  such that

$$x^{(i)} \approx \Phi z^{(i)}.$$

Optimization problem: find matrices  $\Phi, Z$  that minimize

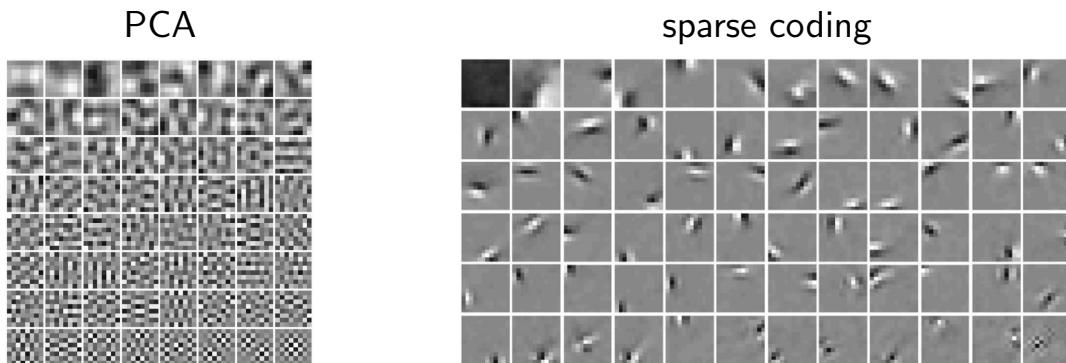
$$\begin{aligned} L(\Phi, Z) &= \|X - \Phi Z\|_F^2 - \lambda \cdot \text{sparsity}(Z) \\ &= \sum_{i=1}^n \left( \|x^{(i)} - \Phi z^{(i)}\|^2 - \lambda \cdot \text{sparsity}(z^{(i)}) \right) \end{aligned}$$

Alternating minimization procedure:

- Initialize  $\Phi$  somehow
- Repeat until convergence:
  - Fixing  $\Phi$ , minimize  $L(\cdot)$  over  $Z$
  - Fixing  $Z$ , minimize  $L(\cdot)$  over  $\Phi$

## Example: image patches

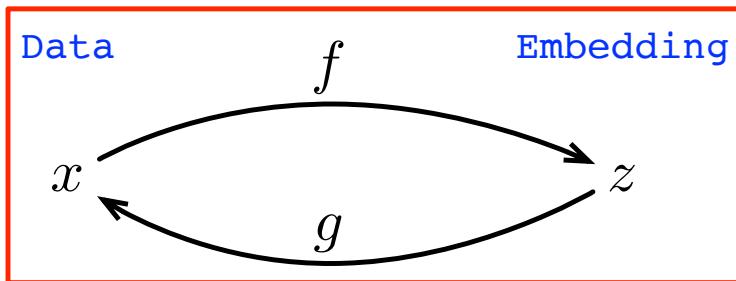
Olshausen-Field (1996), Lewicki-Olshausen (1999): PCA versus sparse coding for natural image patches.



Sparse coding does a much better job at finding a basis that resembles the receptive fields of simple cells in visual cortex.

## C: Variational autoencoder

## Probabilistic autoencoders



- $f$ : encoding map
- $g$ : decoding map

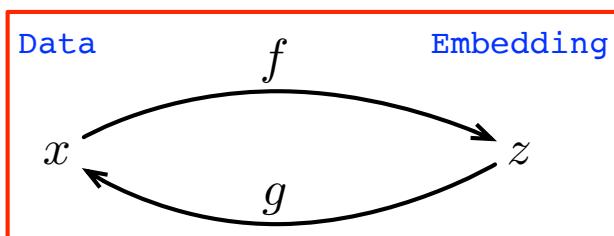
Now  $g(z) \in \Theta$ , where  $\{p_\theta : \theta \in \Theta\}$  is a family of probability distributions over  $x$ -space

Typically: find  $f, g$  to minimize  $\mathbb{E}_X[-\ln p_{g(f(x))}(X)]$

In fact,  $f$  can also be a probabilistic map

## Topic model as probabilistic autoencoder

A document is given by a count-vector  $x \in \mathbb{N}^V$

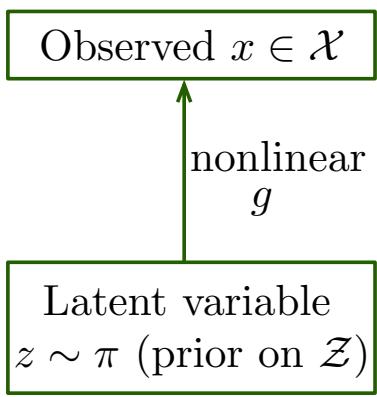


What are  $f$  and  $g$ ?

## Aside: PCA AND $k$ -means

How can we view these as probabilistic autoencoders?

## Nonlinear decoding



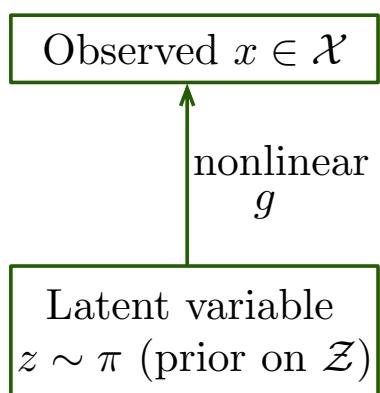
- Let  $\{p_\theta : \theta \in \Theta\}$  be a family of distributions on  $\mathcal{X}$
- Learn nonlinear map  $g : \mathcal{Z} \rightarrow \Theta$  (e.g. neural net)
- $\Pr(x, z|g) = \pi(z)p_{g(z)}(x)$

Example:

- $\mathcal{Z} = \mathbb{R}^k$  and  $\pi$  is Gaussian  $N(0, I_k)$
- $\mathcal{X} = \{0, 1\}^d$
- $\Theta = [0, 1]^d$ ,  $p_\theta$  is product-Bernoulli:

$$p_\theta(x) = \prod_{i=1}^d \theta_i^{x_i} (1 - \theta_i)^{1-x_i}$$

## Approximate encoder



- $\{p_\theta : \theta \in \Theta\}$ , family of distributions on  $\mathcal{X}$
- Neural net  $g : \mathcal{Z} \rightarrow \Theta$
- $\Pr(x, z|g) = \pi(z)p_{g(z)}(x)$

In order to learn  $g$ , need  $\Pr(x|g)$ : intractable!

$$\Pr(x|g) = \int_z \pi(z) p_{g(z)}(x) dz$$

Also intractable:  $\Pr(z|x, g) = \Pr(x, z|g)/\Pr(x|g)$

Learn another neural net to approximate  $\Pr(z|x, g)$

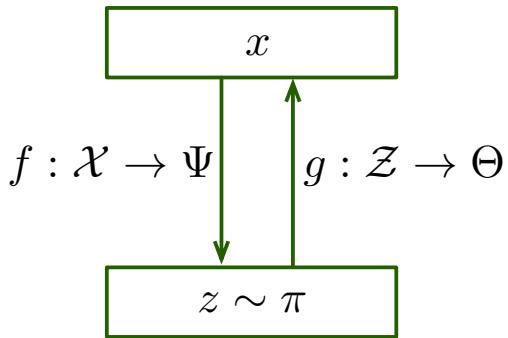
- Let  $\{q_\psi : \psi \in \Psi\}$  be a family of distributions on  $\mathcal{Z}$
- Encoding map  $f : \mathcal{X} \rightarrow \Psi$
- Want:  $q_{f(x)}(z) \approx \Pr(z|x, g)$

E.g. for  $\mathcal{Z} = \mathbb{R}^k$ , could take  $\Psi = \mathbb{R}^{2k}$  with

- $\psi = (\mu_1, \dots, \mu_k, \ln \sigma_1, \dots, \ln \sigma_k)$
- $q_\psi = N(\mu, \text{diag}(\sigma_1^2, \dots, \sigma_k^2))$

## The loss function

Should choose  $g$  to maximize



$$\Pr(x|g), \text{ which equals } \frac{\Pr(x, z|g)}{\Pr(z|x, g)}$$

Instead choose  $g, f$  to maximize

$$L_{f,g}(x) = \mathbb{E}_{z \sim q_{f(x)}} \ln \frac{\Pr(x, z|g)}{q_{f(x)}(z)}$$

This is called the ELBO.

Given data  $x_1, \dots, x_n \in \mathcal{X}$ , solve:

$$\arg \max_{f,g} \sum_{i=1}^n L_{f,g}(x_i)$$

## Property of the ELBO

**Claim:**  $L_{f,g}(x) = \ln \Pr(x|g) - K(q_{f(x)}, \Pr(\cdot|x, g))$

# VAE example: MNIST

8617814828	8165167672	2838385738	8208923700
9683960319	8594682168	8382793538	7519117144
3391368179	8153288433	8599239516	8762082829
8908691963	2868910041	1988833497	2986387461
8233331336	5172018359	2736430263	5479899910
6998616663	6561491758	5970582845	6824988281
9526651899	1343983470	6943628557	7582461383
9989312823	4582970957	8490507066	9939299390
0461232088	6944972393	7436203601	4524390184
9759939851	2645609798	2180971000	8872314236

(a) 2-D latent space      (b) 5-D latent space      (c) 10-D latent space      (d) 20-D latent space

Kingma and Welling (2014). Autoencoding variational Bayes.

# VAE example: MNIST

Kingma and Welling (2014). Autoencoding variational Bayes.