

Autoencoders

Autoencoders are neural network that are trained to learn how to map their input to their input. Internally, it has an hidden layer h that contains a lossy summary of the relevant feature for the task.

An autoencoder can be seen has a two parts network

- Encoder function: $h = f(x)$
- Decoder function: $\tilde{x} = g(h)$

The simplest autoencoder is a MLP:

$$\begin{aligned}h &= \sigma_1(W_{xh}x) \\ \tilde{x} &= \sigma_2(W_{hx}h)\end{aligned}\tag{1}$$

Autoencoders

- Encoder function: $\mathbf{h} = f(\mathbf{x})$
- Decoder function: $\tilde{\mathbf{x}} = g(\mathbf{h})$

The simplest autoencoder is a MLP:

$$\begin{aligned}\mathbf{h} &= \sigma_1(W_{xh}\mathbf{x}) \\ \tilde{\mathbf{x}} &= \sigma_2(W_{hx}\mathbf{h})\end{aligned}\tag{2}$$