

Which Inputs of a Neural Net Lead to the Same Output?

Raymond Veldhuis, Zohra Rezgui, Amina Bassit

April 14, 2021

1 Introduction

In this assignment, we want to analyse what type of variations of the input of a neural network lead to the same output. This is important, because sometimes this behaviour is desired, but at other times, it is not because it is unpredictable. An example of a desired mapping of input variations to the same output is robustness to variations of pose or illumination of a facial image in face recognition. This is a desired behaviour, because we don't want the pose or the illumination to affect the recognition result.

2 Background

Neural networks that operate on images, such as autoencoders and classification networks map input images, denoted as vectors $\mathbf{x} \in \mathbb{R}^D$ onto a lower dimensional *latent vector*, sometimes also called *embedding*, $\mathbf{z} = F(\mathbf{x}) \in \mathbb{R}^d$, with $d < D$ and often $d \ll D$. The latent vector is then further used for, for instance, reconstruction or classification. Ideally, the mapping is such that in \mathbf{z} the relevant properties of \mathbf{x} for the task at hand (e.g. reconstruction or classification) are preserved. The network function $F(\mathbf{x})$ is noninvertible and more than one input \mathbf{x} can be mapped onto the same latent variable \mathbf{z} as shown in Figure 1. Each row of Figure 2 shows the variation of facial images that could lead to the same or very close latent vectors as an artificially generated example.

A bit of math: For a certain \mathbf{z}_0 we denote the inputs that are mapped onto \mathbf{z} by $F^{-1}\{\mathbf{z}_0\} = \{\mathbf{x} | F(\mathbf{x}) = \mathbf{z}_0\}$. Note that F^{-1} is not the inverse of F , because F is not invertible. The function F^{-1} operates on sets and is called the *preimage* of a set in mathematics. Formally, $F^{-1}A = \{\mathbf{x} | F(\mathbf{x}) \in A\}$. Here we are interested in $F^{-1}\{\mathbf{z}_0\}$ and maybe, if time permits in $F^{-1}\{\mathbf{z} | \|\mathbf{z} - \mathbf{z}_0\| \leq \epsilon\}$: the set of inputs mapped onto a small neighbourhood of \mathbf{z}_0 .

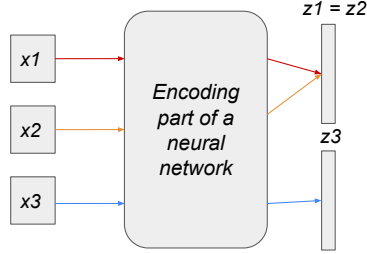


Figure 1: Illustration of samples x_1 and x_2 mapped to the same latent vector by a neural network.

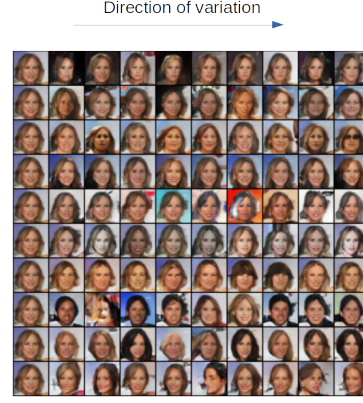


Figure 2: An artificially generated example of facial images that could be mapped on to the same latent variable. Every row corresponds to one latent variable.

3 The assignment

The goal of this assignment is to study preimages of elements of the latent space of a classifier, and, if time permits, preimages of small neighbourhoods of elements of the latent space of a classifier, by visualisation of these variations in the input image. So, in human language: *We want to see which images are mapped onto the same latent variable.*

More specifically, for a relatively simple classifier, i.e. with a low-dimensional latent space, and for a manageable dataset, in this case the MNIST set of digits, the task is to analyse the preimages of the latent vectors corresponding to the input digits by visualisation.

We suggest to do this by linearisation of the network function F . If we have $F(\mathbf{x}_0) = \mathbf{z}_0$, then for small deviations Δ of \mathbf{x}_0 we may write

$$F(\mathbf{x}_0 + \Delta) \simeq \mathbf{z}_0 + \frac{dF(\mathbf{x}_0)}{d\mathbf{x}} \Delta.$$

The first derivative $d \times D$ -matrix $\frac{dF(\mathbf{x}_0)}{d\mathbf{x}}$ can be found analytically by applying back propagation to the network. We look for vectors Δ that are in the null-space of $\frac{dF(\mathbf{x}_0)}{d\mathbf{x}}$, because for those Δ we have, by approximation, that $\mathbf{x}_0 + \Delta \in F^{-1}\{\mathbf{z}_0\}$.

More details can be provided when the research proposal is written.