

AutoEncoders for Step Functions

IA311 Project

Hélène Maxcici

Télécom Paris

February 27, 2022

Overview

1. Introduction

- 1.1 Autoencoders
- 1.2 Goal
- 1.3 One Layer Encoder

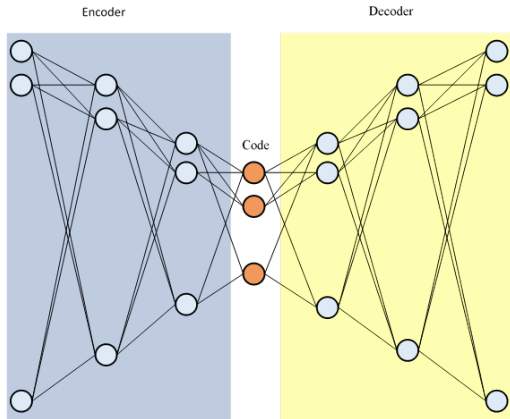
2. One Layer Decoder

3. Two Layers Decoder

4. Conclusion

Autoencoders

- Unsupervised learning of data compression
- Success in image denoising, dimensionality reduction, anomaly detection...
- However, the code is not always understandable and linear \implies not manipulable



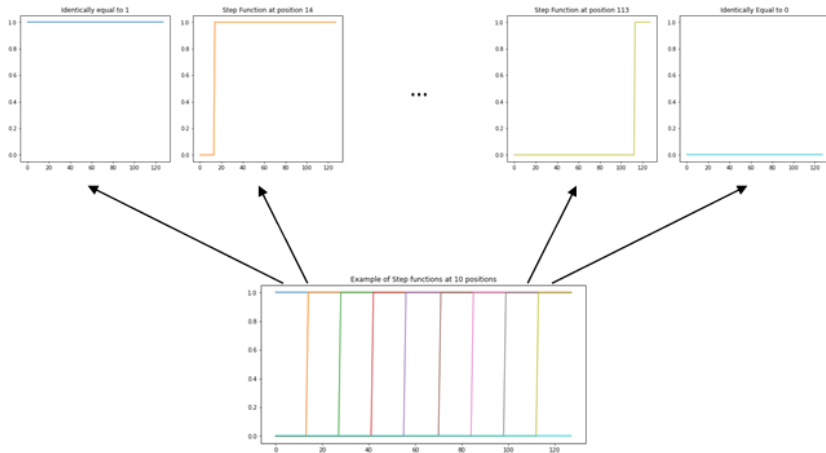
Goal

The goal is to build autoencoder for step functions with one varying parameter s.t.

- The reconstruction is perfect
- The autoencoder is the simplest it can be
- The code is scalar and linear with the parameter
- The output activation is a ReLU

Step Functions

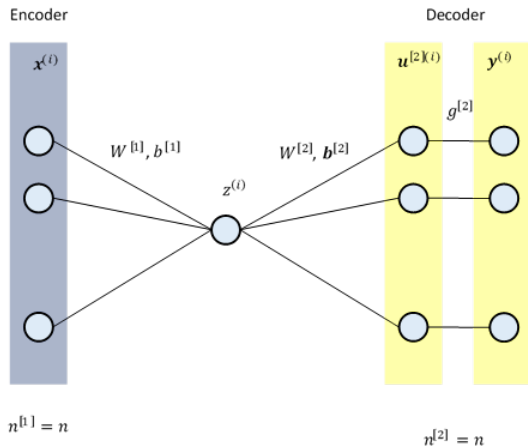
We will consider $n + 1$ step functions of dimension n and same amplitude C .



One Layer Encoder

- A one MLP layer encoder is enough to generate linear codes with respect to the positions
- A non-linearity is not required in the encoder

One Layer Decoder



Necessary Condition

Lemma

A one layer decoder is able to generate the $n + 1$ different step functions out of their scalar codes z only if there exist two non-overlapping proper intervals A and B such that

$$g^{[2]}(u) = \begin{cases} 0, & \text{if } u \in A \\ C, & \text{if } u \in B \end{cases}$$

Sigmoid Activation Function

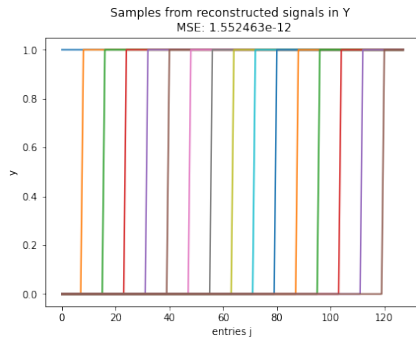
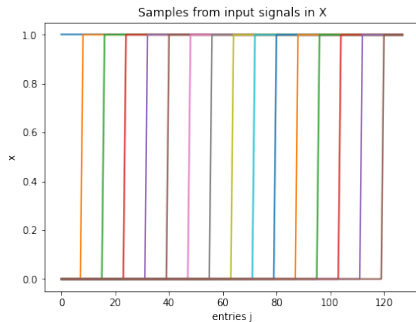
Proposition

Let the output activation be a *sigmoid* defined by $g^{[2]}(u) = \frac{e^u}{1+e^u}$. And let the encoding be linear. Then, we can find a solution $W^{[2]}, \mathbf{b}^{[2]} \in \mathbb{R}^n$ for the decoder, for which there exists an $\epsilon \in \mathbb{R}^+$ arbitrarily small such that

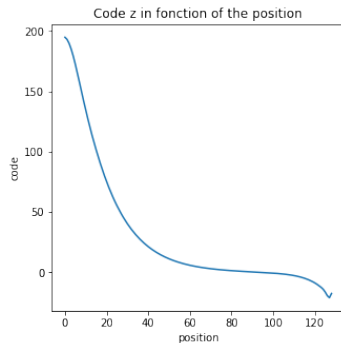
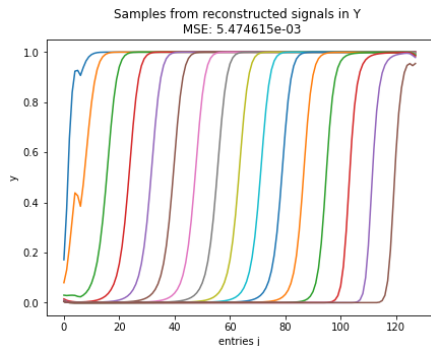
$$MSE = \frac{1}{n(n+1)} \sum_{i=1}^{n+1} \|\mathbf{y}^{(i)} - \mathbf{x}^{(i)}\|_2^2 < \epsilon$$

Handcrafted Solution

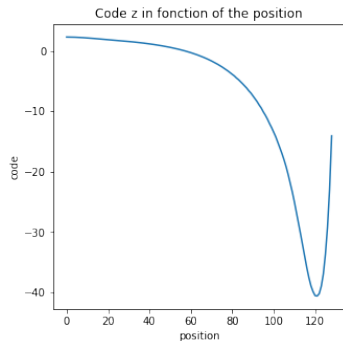
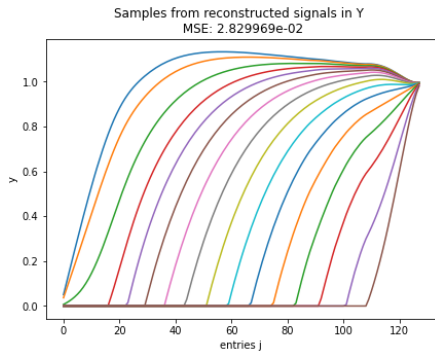
- Linear encoder : $z^{(i)} = i$
- MSE upper bound : $\epsilon = 10^{-5}$



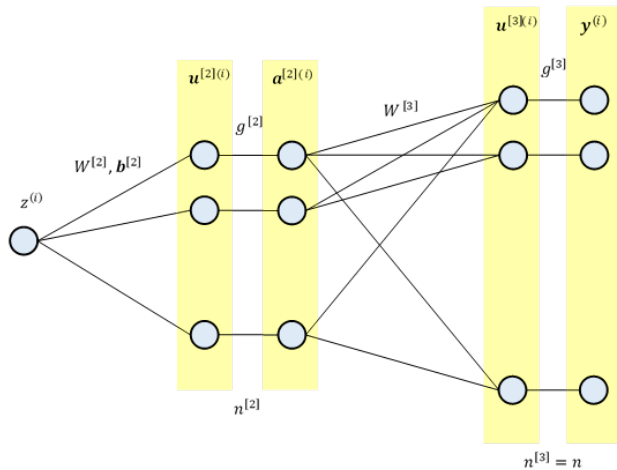
Training with no constraints



Training with no constraints



Two Layers Decoder



Layer 3

- Perfect Reconstruction:

$$Y_{(n+1) \times n} = X_{(n+1) \times n} = \begin{bmatrix} C & C & \dots & C & C \\ 0 & C & \dots & C & C \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 0 & C \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix}$$

- ReLU Activation Function:

$$Y_{(n+1) \times n} = \text{ReLU}(U_{(n+1) \times n}^{[3]}) = \begin{bmatrix} \max(0, \mathbf{u}_1^{[3](1)}) & \max(0, \mathbf{u}_2^{[3](1)}) & \dots & \max(0, \mathbf{u}_n^{[3](1)}) \\ \max(0, \mathbf{u}_1^{[3](2)}) & \max(0, \mathbf{u}_2^{[3](2)}) & \dots & \max(0, \mathbf{u}_n^{[3](2)}) \\ \vdots & \vdots & \ddots & \vdots \\ \max(0, \mathbf{u}_1^{[3](n)}) & \max(0, \mathbf{u}_2^{[3](n)}) & \dots & \max(0, \mathbf{u}_n^{[3](n)}) \\ \max(0, \mathbf{u}_1^{[3](n+1)}) & \max(0, \mathbf{u}_2^{[3](n+1)}) & \dots & \max(0, \mathbf{u}_n^{[3](n+1)}) \end{bmatrix}$$

Layer 3

Lemma

A perfect reconstruction $Y = X$, is achieved only if $U^{[3]}$ is **full rank** matrix of the following form,

$$U^{[3]} = \begin{bmatrix} C & C & \dots & C & C \\ \mathbf{u}_1^{[3](2)} & C & \dots & C & C \\ \mathbf{u}_1^{3} & \mathbf{u}_2^{3} & \dots & C & C \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{u}_1^{[3](n)} & \mathbf{u}_2^{[3](n)} & \dots & \mathbf{u}_{n-1}^{[3](n)} & C \\ \mathbf{u}_1^{[3](n+1)} & \mathbf{u}_2^{[3](n+1)} & \dots & \mathbf{u}_{n-1}^{[3](n+1)} & \mathbf{u}_n^{[3](n+1)} \end{bmatrix} \quad (1)$$

where, $\mathbf{u}_j^{[3](i)} \leq 0$ for all $i = 2, \dots, n+1$ and $j = 1, \dots, i-1$.

Layer 3

- Layer 3 linear projection:

$$U^{[3]} = A^{[2]} W^{[3]T}$$

- Rank of a matrix product:

$$\text{rank}(U^{[3]}) \leq \min \left(\text{rank} \left(A^{[2]}_{(n+1) \times n^{[2]}} \right), \text{rank} \left(W^{[3]}_{n \times n^{[2]}} \right) \right)$$

Lemma

The preactivation matrix $U^{[3]}$ can be full rank under three necessary conditions:

1. $n^{[2]} = n$
2. $W^{[3]}$ is full rank
3. $A^{[2]}$ is full rank

Layer 2

- Layer 2:

$$A^{[2]} = g^{[2]} \left(\begin{bmatrix} z^{(1)} \\ z^{(2)} \\ \vdots \\ z^{(n+1)} \end{bmatrix} W^{[2]T} + \begin{bmatrix} \dots & \mathbf{b}^{[2]} & \dots \\ \dots & \mathbf{b}^{[2]} & \dots \\ & \vdots & \\ \dots & \mathbf{b}^{[2]} & \dots \end{bmatrix} \right)$$

Lemma

If $g^{[2]} = \text{ReLU}$, then the presence of a bias is a necessary and sufficient condition to make $A^{[2]}$ full rank.

Layer 2

•

$$A^{[2]} = \begin{bmatrix} g^{[2]}(W_1^{[2]}(z^{(1)} + k_1)) & g^{[2]}(W_2^{[2]}(z^{(1)} + k_2)) & \dots & g^{[2]}(W_n^{[2]}(z^{(1)} + k_n)) \\ g^{[2]}(W_1^{[2]}(z^{(2)} + k_1)) & g^{[2]}(W_2^{[2]}(z^{(2)} + k_2)) & \dots & g^{[2]}(W_n^{[2]}(z^{(2)} + k_n)) \\ \vdots & \vdots & & \vdots \\ g^{[2]}(W_1^{[2]}(z^{(n+1)} + k_1)) & g^{[2]}(W_2^{[2]}(z^{(n+1)} + k_2)) & \dots & g^{[2]}(W_n^{[2]}(z^{(n+1)} + k_n)) \end{bmatrix}$$

where $k_j = \frac{\mathbf{b}_j^{[2]}}{W_j^{[2]}}$ for $j \in \{1, \dots, n\}$.

Theoretical Conclusion on Architecture

Proposition

Let's drop the signal identically equal to 0 from the dataset. A two layers decoder can perfectly reconstruction step functions from their scalar codes if the layer 2 satisfies the following conditions:

1. Layer 2 has at least $n^{[2]} = n$ neurons, where n is the dimension of the signal.
2. The activation function $g^{[2]}$ of layer 2 is non-linear. In particular, $g^{[2]} = \text{ReLU}$.
3. Layer 2 has a bias vector.

Bias in Layer 3

- A bias in Layer three is not necessary. It will only increase the complexity.

Handcrafting the solutions

- Layer 3 linear projection:

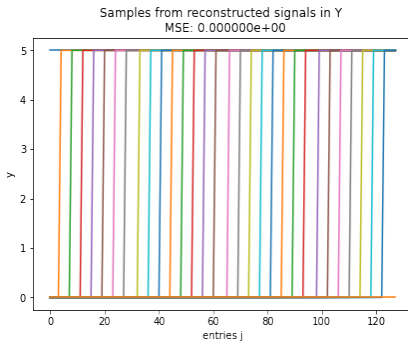
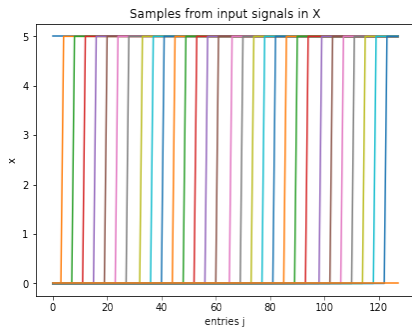
$$U^{[3]} = A^{[2]} W^{[3]T}$$

- Omitting the signal identically equal to 0,

$$W^{[3]T} = A^{[2]-1} U^{[3]}$$

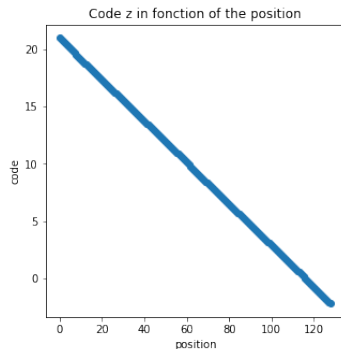
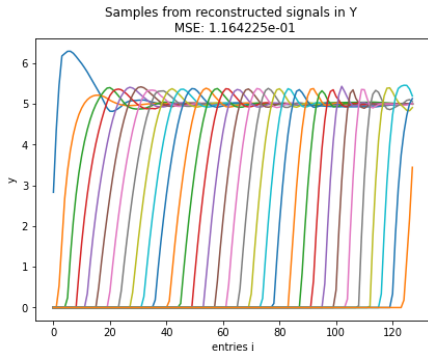
Handcrafted Solution - Linear Encoder

- Linear encoder : $z^{(i)} = i$
- $W^{[2]}$ and $\mathbf{b}^{[2]}$ makes $A^{[2]}$ full rank
- $U^{[3]}$ arbitrary of the form (1)



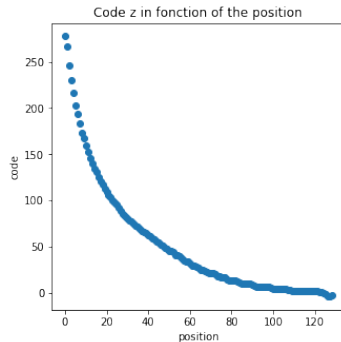
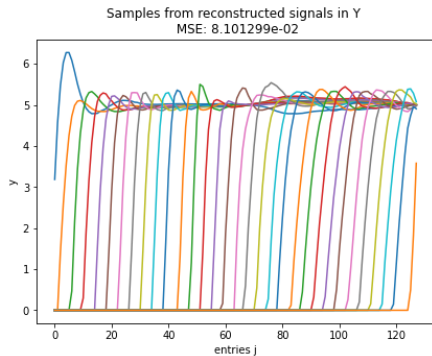
Trained AE - Constrained Linear Encoder

- Linear encoder : $z^{(i)} = \alpha i + \beta$
- RMSProp optimizer with a learning rate of 10^{-5}
- Number of epochs: 2500000



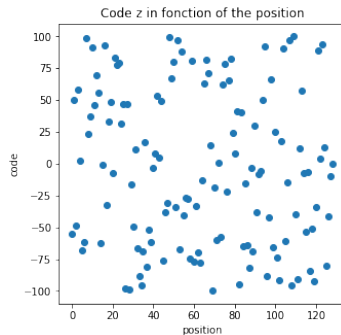
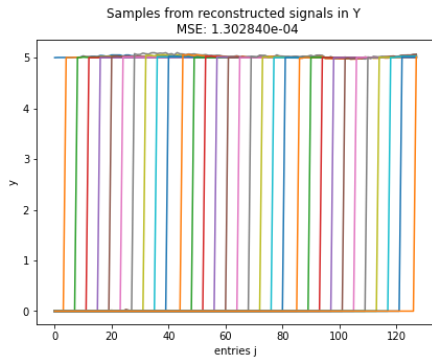
Trained AE - No constraints

- RMSProp optimizer with a learning rate of 10^{-4}
- Number of epochs: 2500000



Handcrafted Solution - Random Encoder

- Random Encoder
- Sort $z^{(i)}$ then find a solution



Set of solutions

Conjecture

$A^{[2]}$ is full rank if and only if at most two rows and at most two columns have zeros at identical positions and the ratios $\left\{ \frac{\mathbf{b}_j^{[2]}}{W_j^{[2]}} \right\}_{j=1,\dots,n}$ are mutually different at the non zero entries.

$$A^{[2]} = \begin{bmatrix} g^{[2]}(W_1^{[2]}(z^{(1)} + k_1)) & g^{[2]}(W_2^{[2]}(z^{(1)} + k_2)) & \dots & g^{[2]}(W_n^{[2]}(z^{(1)} + k_n)) \\ g^{[2]}(W_1^{[2]}(z^{(2)} + k_1)) & g^{[2]}(W_2^{[2]}(z^{(2)} + k_2)) & \dots & g^{[2]}(W_n^{[2]}(z^{(2)} + k_n)) \\ \vdots & \vdots & & \vdots \\ g^{[2]}(W_1^{[2]}(z^{(n+1)} + k_1)) & g^{[2]}(W_2^{[2]}(z^{(n+1)} + k_2)) & \dots & g^{[2]}(W_n^{[2]}(z^{(n+1)} + k_n)) \end{bmatrix}$$

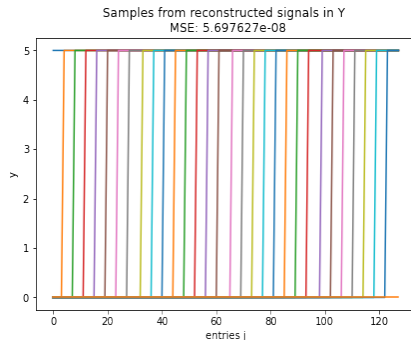
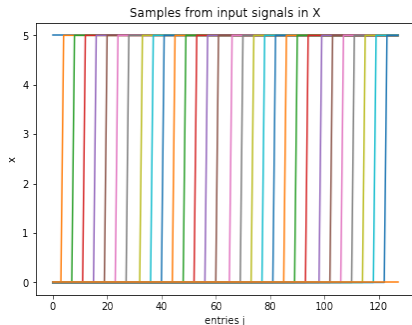
Set of solutions

Lemma

If $(W^{[2]}, \mathbf{b}^{[2]}, W^{[3]})$ is a solution for a decoder with a perfect reconstruction, then for all positive diagonal matrices D , $(DW^{[2]}, \mathbf{b}^{[2]}D, W^{[3]}D^{-1})$ is also a solution.

Generated Solutions

- $(W^{[2]}, \mathbf{b}^{[2]}, W^{[3]})$ from handcrafted solution with linear encoder
- D random positive diagonal matrix



Set of solutions

- $\mathbf{k} = [k_1, k_2, \dots, k_n] = \left[\frac{\mathbf{b}_1^{[2]}}{W_1^{[2]}}, \frac{\mathbf{b}_2^{[2]}}{W_2^{[2]}}, \dots, \frac{\mathbf{b}_n^{[2]}}{W_n^{[2]}} \right]^T$

Corollary

Let \mathcal{U} be the set of all $U^{[3]}$ having the form (1). And let \mathcal{K} be the set of all \mathbf{k} that make the activation matrix of layer 2 full rank. For each fixed couple $(\mathbf{k}, U^{[3]}) \in \mathcal{K} \times \mathcal{U}$, there exists a set of infinitely many solutions described by the two following equations,

$$W^{[2]} = W^{[3]-1} W^{[2]*}$$

$$\mathbf{b}^{[2]} = \mathbf{b}^{[2]*} W^{[3]-1 T}$$

where $W^{[2]*}$ is a constant matrix of dimension $n \times 1$ and $\mathbf{b}^{[2]*}$ is a constant vector of dimension n .

Set of solutions

•

$$A^{[2]} = \begin{bmatrix} g^{[2]}(W_1^{[2]}(z^{(1)} + k_1)) & g^{[2]}(W_2^{[2]}(z^{(1)} + k_2)) & \dots & g^{[2]}(W_n^{[2]}(z^{(1)} + k_n)) \\ g^{[2]}(W_1^{[2]}(z^{(2)} + k_1)) & g^{[2]}(W_2^{[2]}(z^{(2)} + k_2)) & \dots & g^{[2]}(W_n^{[2]}(z^{(2)} + k_n)) \\ \vdots & \vdots & & \vdots \\ g^{[2]}(W_1^{[2]}(z^{(n+1)} + k_1)) & g^{[2]}(W_2^{[2]}(z^{(n+1)} + k_2)) & \dots & g^{[2]}(W_n^{[2]}(z^{(n+1)} + k_n)) \end{bmatrix}$$

•

$$W^{[3]T} = \left(I - \frac{1}{1 + \text{trace}(\Delta W^{[3]+T} U^{[3]-1})} W^{[3]+T} U^{[3]-1} \Delta \right) W^{[3]+T}$$

Conclusion

- Three layers Autoencoder: simplest + admissible
- Infinitely many solutions \implies Too complex?
- Very slow training
- High tendency for a nonlinear encoder
- No requirements on the encoding

Thanks For Your Attention!

Any Questions?