

Autoencoders

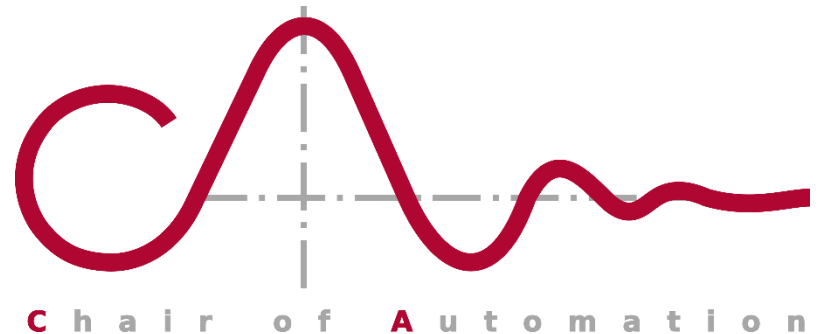
An Introduction

Author: Paul O'Leary and Anika Terbuch
Date: 6 Jule 2022
Document Number: CoA2020-MLSS
Conference: Summer School Presentation

Chair of Automation

Department Product Engineering, University of Leoben
Peter Tunner Straße 27, 8700 Leoben, Austria

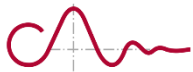
phone: +43/(0)3842/402-5301
fax: +43/(0)3842/402-5302
email: automation@unileoben.ac.at
web: automation.unileoben.ac.at





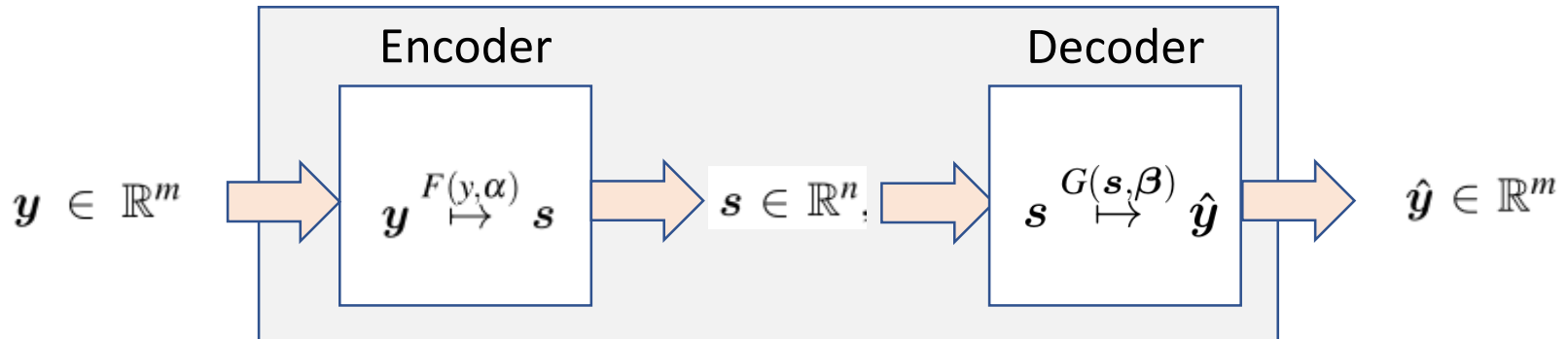
Context

1. The primary focus in this presentation is on time-series analysis.
2. Application to real-time machine data that for multi-variate time-series (MVTs).
3. Autoencoders will be considered as a means of identifying relevant information in data and with this to enable dimensionality reduction.
4. Autoencoders can be trained in an unsupervised manner, this alleviates the need for labelled data.
5. Final goal is hybrid-learning, i.e., combining: a) a-priori knowledge, b) analytical techniques and c) machine learning.



Elements of an Autoencoder

Autoencoder structure

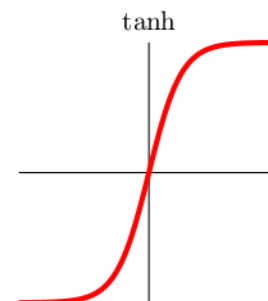
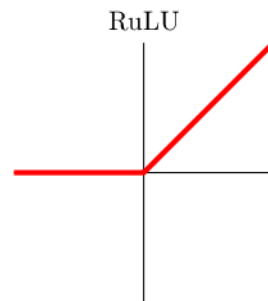


Single layer model

$$\mathbf{s} = f(\mathbf{W}_e \mathbf{y} + \mathbf{b}_e) \in \mathbb{R}^n$$

$$\hat{\mathbf{y}} = g(\mathbf{W}_d \mathbf{y} + \mathbf{b}_d) \in \mathbb{R}^m$$

Both $f(x)$ and $g(x)$ are activation functions, e.g.



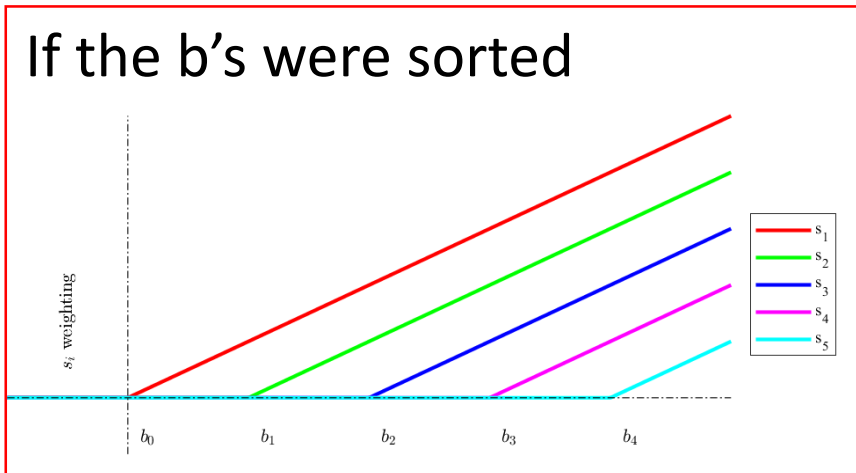
Note both these activation functions have significant regions which are linear

Intuitive approach: Piece Wise Sum of Functions

With the ReUL activation we have

$$z_i = \mathbf{w}_i + b_i$$
$$s_i = f(z_i)$$
$$\Rightarrow s_i = \mathbf{w}_i \mathbf{y} + b_i \quad \text{if } s_i > 0.$$

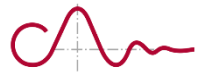
If the b 's were sorted



If the b 's were random



With the ReUL we always have a piece wise function of linear operations on \mathbf{y} .



First use in the Context of ML

First definition in a learning context:

„ Learning Internal Representations by Error Propagation“ by:
Rumelhart, Hinton and Williams, 1965.

TABLE 5

Input Patterns		Hidden Unit Patterns		Output Patterns
1000000	->	.5 0 0	->	1000000
0100000	->	0 1 0	->	0100000
0010000	->	1 1 0	->	0010000
0001000	->	1 1 1	->	0001000
0000100	->	0 1 1	->	0000100
0000010	->	.5 0 1	->	0000010
0000001	->	1 0 .5	->	0000001
0000000	->	0 0 .5	->	0000000

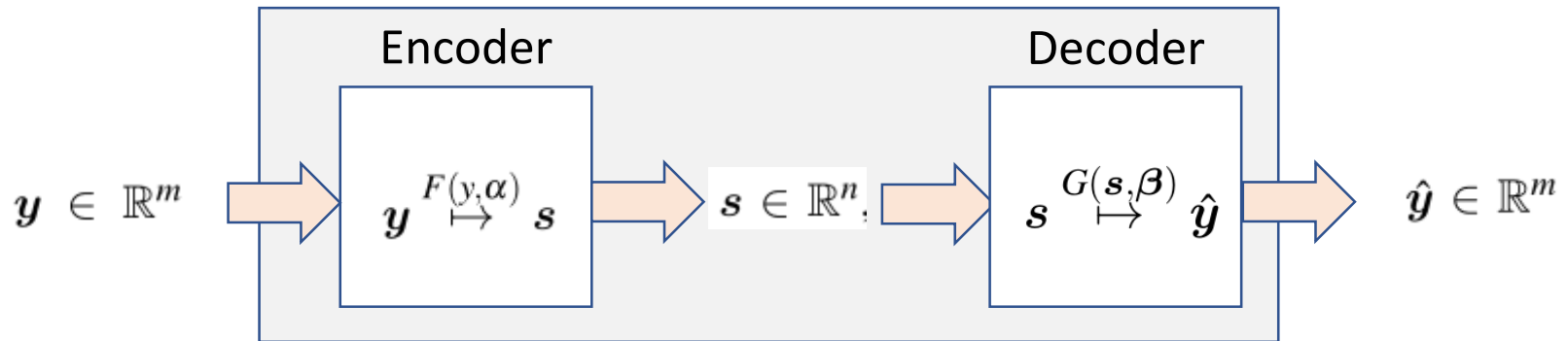
TABLE 8

Input Patterns		Hidden Unit Patterns		Output Patterns
00 + 00	-	111	-	000
00 + 01	-	110	-	001
00 + 10	-	011	-	010
00 + 11	-	010	-	011
01 + 00	-	110	-	001
01 + 01	-	010	-	010
01 + 10	-	010	-	011
01 + 11	-	000	-	100
10 + 00	-	011	-	010
10 + 01	-	010	-	011
10 + 10	-	001	-	100
10 + 11	-	000	-	101
11 + 00	-	010	-	011
11 + 01	-	000	-	100
11 + 10	-	000	-	101
11 + 11	-	000	-	110

@incollection{rumelhart:errorpropnote, address = {Cambridge, MA}, author = {Rumelhart, David E. and Hinton, Geoffrey E. and Williams, Ronald J.}, booktitle = {Parallel Distributed Processing: Explorations in the Microstructure of Cognition, {V}olume 1: {F}oundations}, editor = {Rumelhart, David E. and McClelland, James L.}, pages = {318--362}, year = 1985, publisher = {MIT Press}, title = {Learning Internal Representations by Error Propagation},}

Elements of an Autoencoder

Autoencoder structure (mappings)



Composite map

$$\hat{y} = G(F(y, \alpha), \beta)$$

The goal

The goal is to generate a reproduction \hat{y} of the input y that achieves a high degree of dimensionality reduction, i.e., $n < m$; while, maintaining the *significant information* from the input data.

Expressing the goal as a cost function

$$E(\alpha, \beta) = \|\mathbf{r}\|_2 + \lambda R(\mathbf{s}, N(0, 1))$$

The Concept of Low Rank Identity

Can the identity matrix I be factored?

Of course, since $B^{-1} B = I$. With special case $A^T A = I$.

Do low rank factorizations exist?

Yes, if we can accept some finite error ϵ .

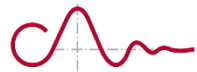
$$\begin{array}{c} n \\ \boxed{A} \\ m \end{array} \begin{array}{c} m \\ \boxed{A^T} \\ n \end{array} = \begin{array}{c} m \\ \boxed{I + \epsilon} \\ m \end{array}$$

TABLE 5

Input Patterns		Hidden Unit Patterns		Output Patterns
1000000	-	.5 0 0	-	1000000
0100000	-	0 1 0	-	0100000
0010000	-	1 1 0	-	0010000
0001000	-	1 1 1	-	0001000
0000100	-	0 1 1	-	0000100
0000010	-	.5 0 1	-	0000010
0000001	-	1 0 .5	-	0000001
0000000	-	0 0 .5	-	0000000

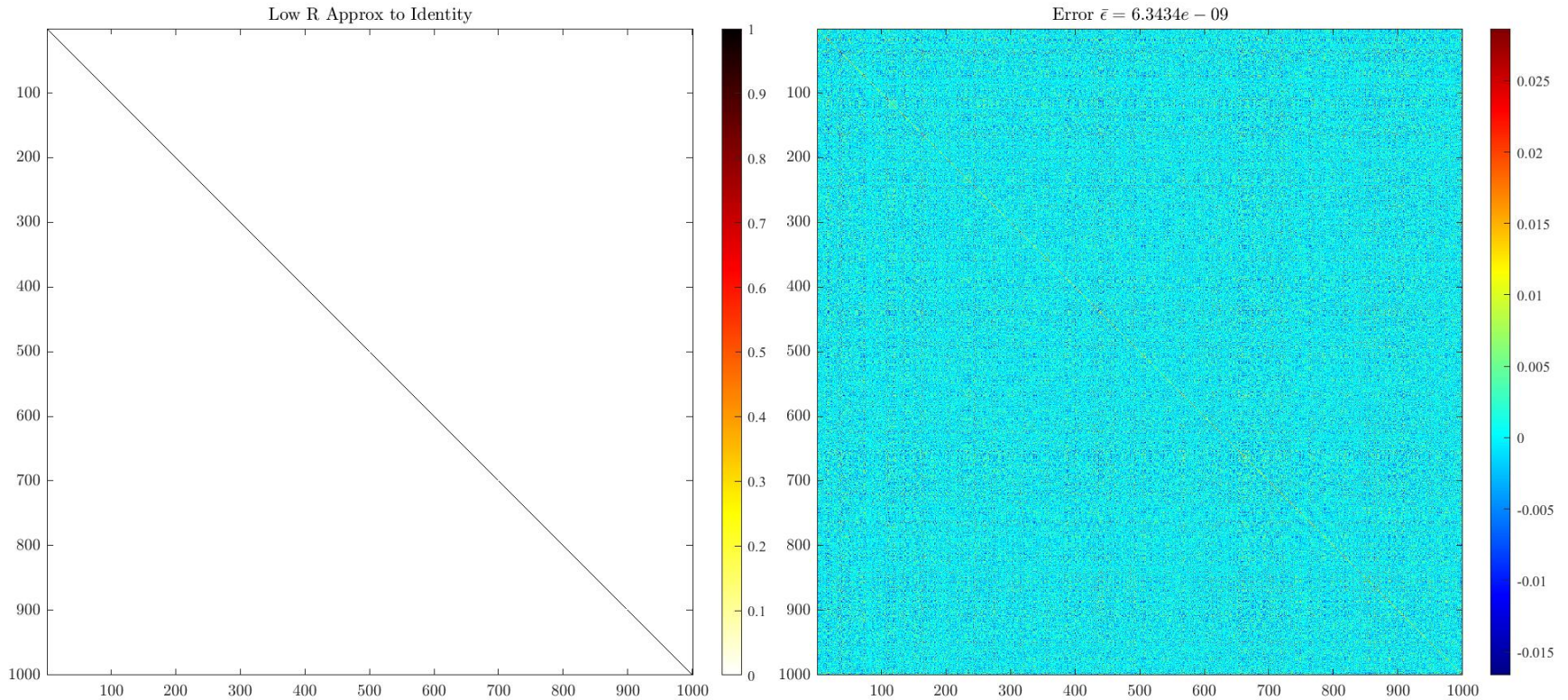
The most efficient low rank approximations are Sylvester equations of the form:

$$\mathbf{A} \mathbf{X} + \mathbf{X} \mathbf{B} = \mathbf{C}$$



Example Low Rank Approximation to Identity

Rank 10 deficient with $\epsilon \approx 2.2e - 9$



Low Rank Approximations are Nonunique

Original

Original



Both reconstructions
are rank 10 deficient

Low rank approx (1)



Random bases

Low rank approx (2)



Polynomial bases

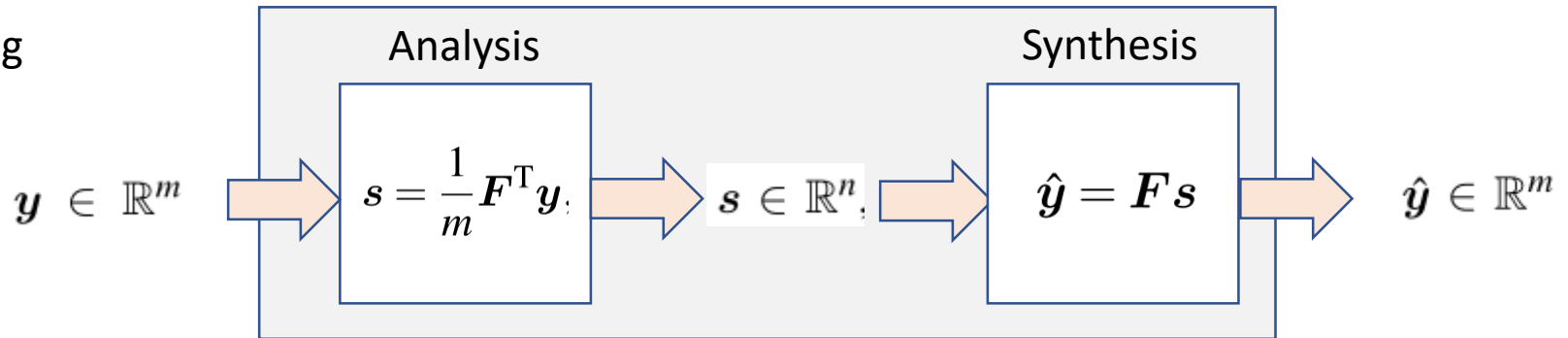
Structure (information)
of the error is important,
not only the magnitude.

Discrete bases functions
can be used to define
structure.

Analogy: Analysis and Synthesis Functions

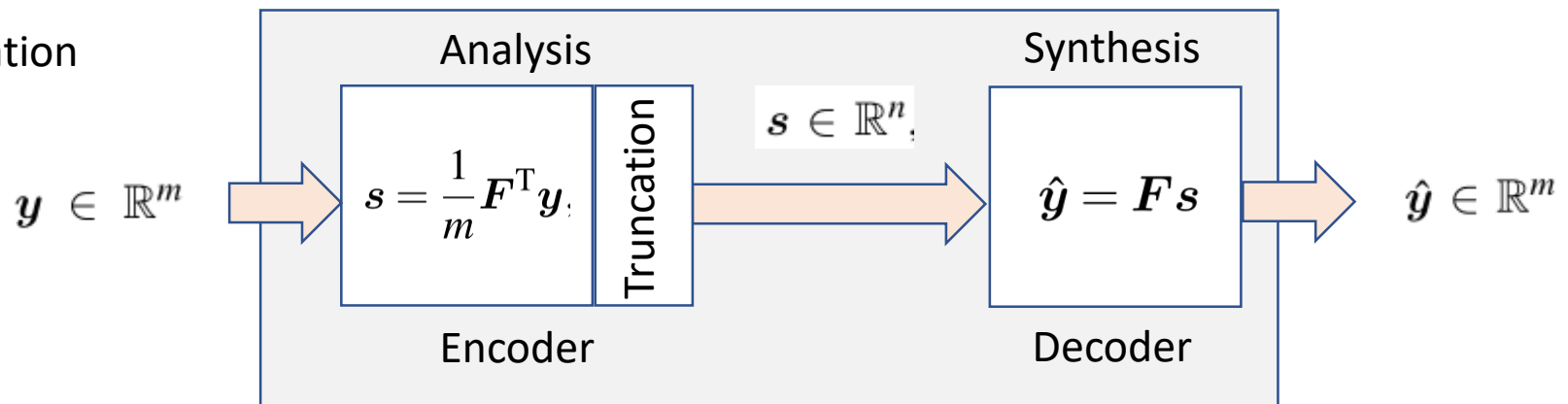
Fourier analysis in the context of autoencoders

Training



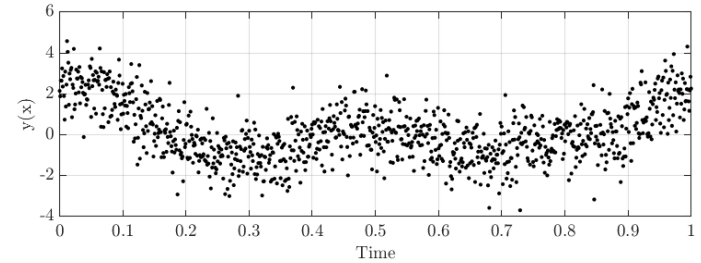
Fourier filtering using truncation (Dimensionality reduction)

Evaluation

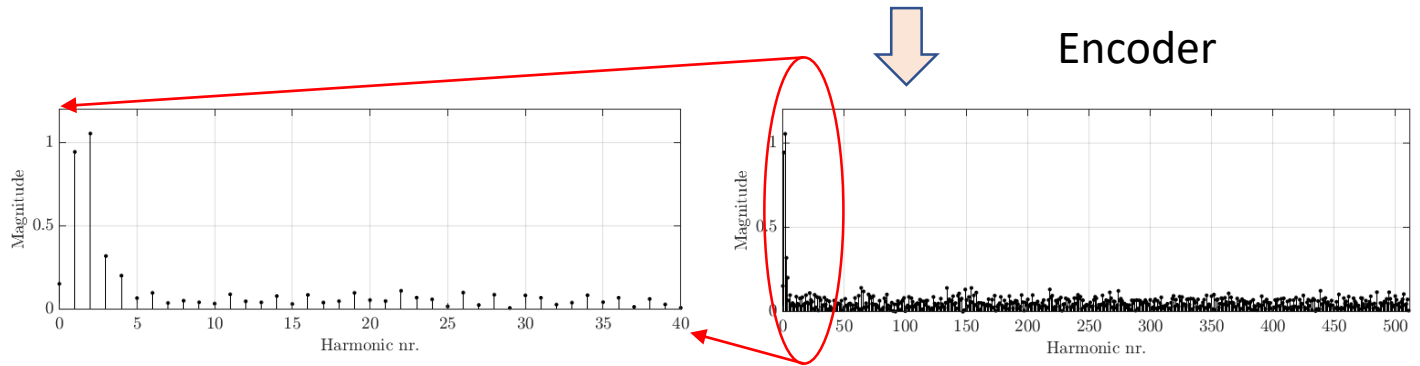


Example of Truncated Fourier Spectrum

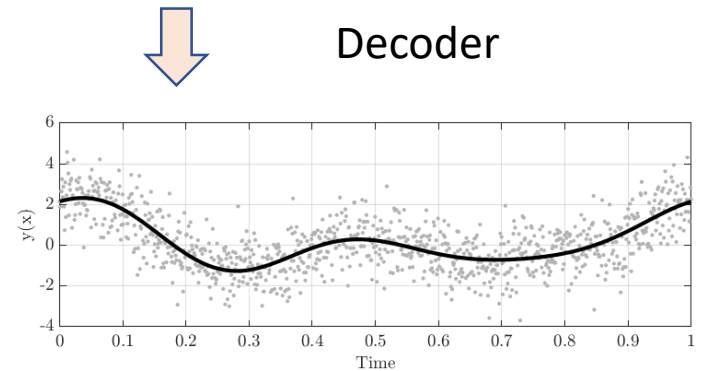
Raw data



Spectrum

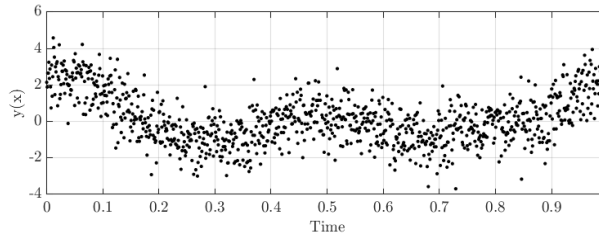


Generation (truncated spectrum)



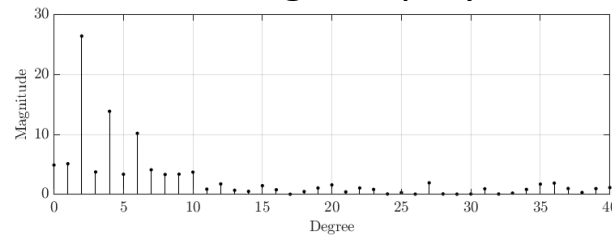
Low Rank Approximations with Bases

Raw data

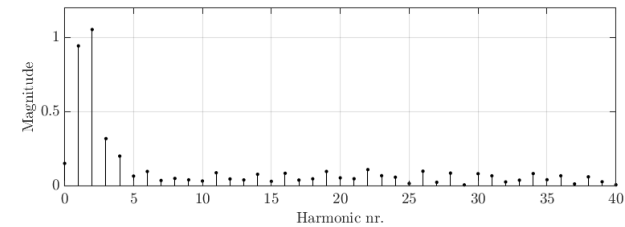


Spectrum wrt Bases

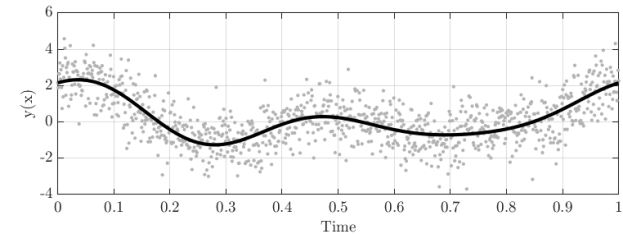
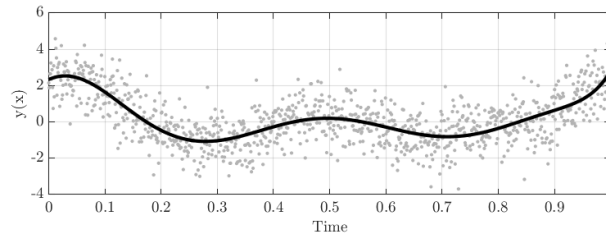
Discrete orthogonal polynomials



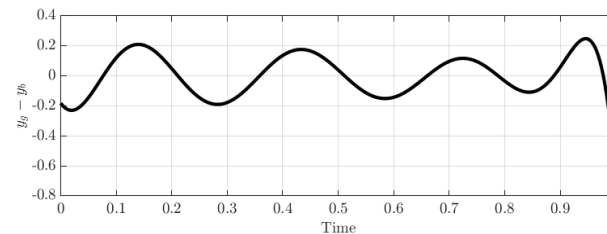
Fourier bases (DFT)



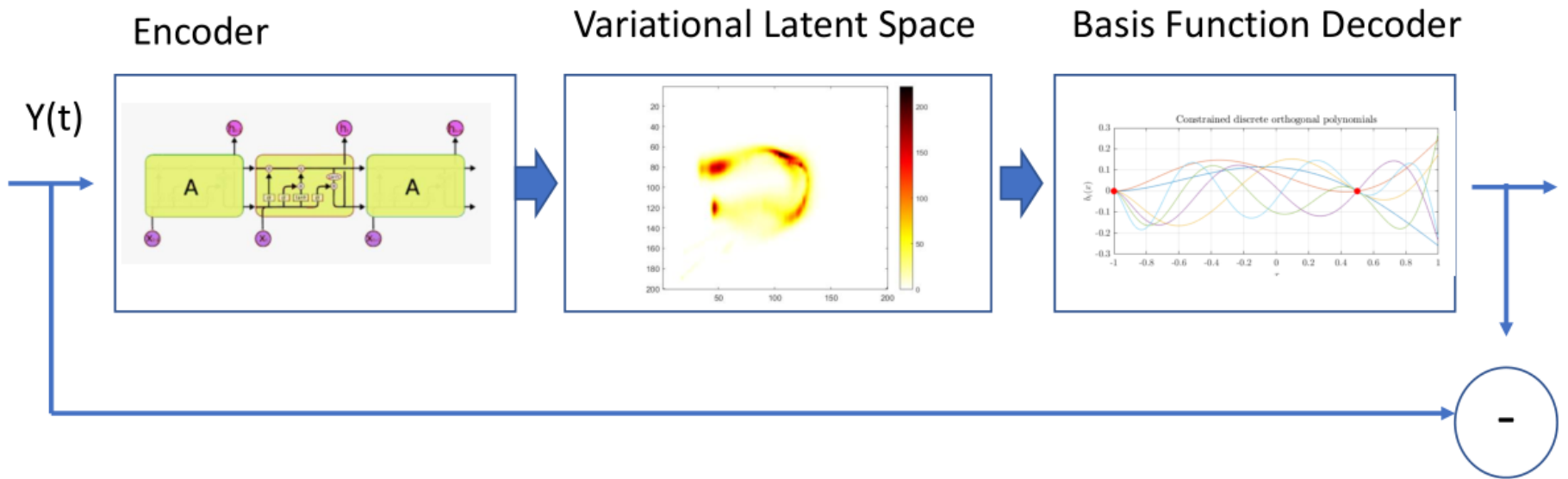
Low rank approximations



Model differences



New Proposed Architecture



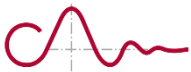
But also $Y(x)$

Error prorogation

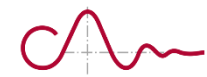
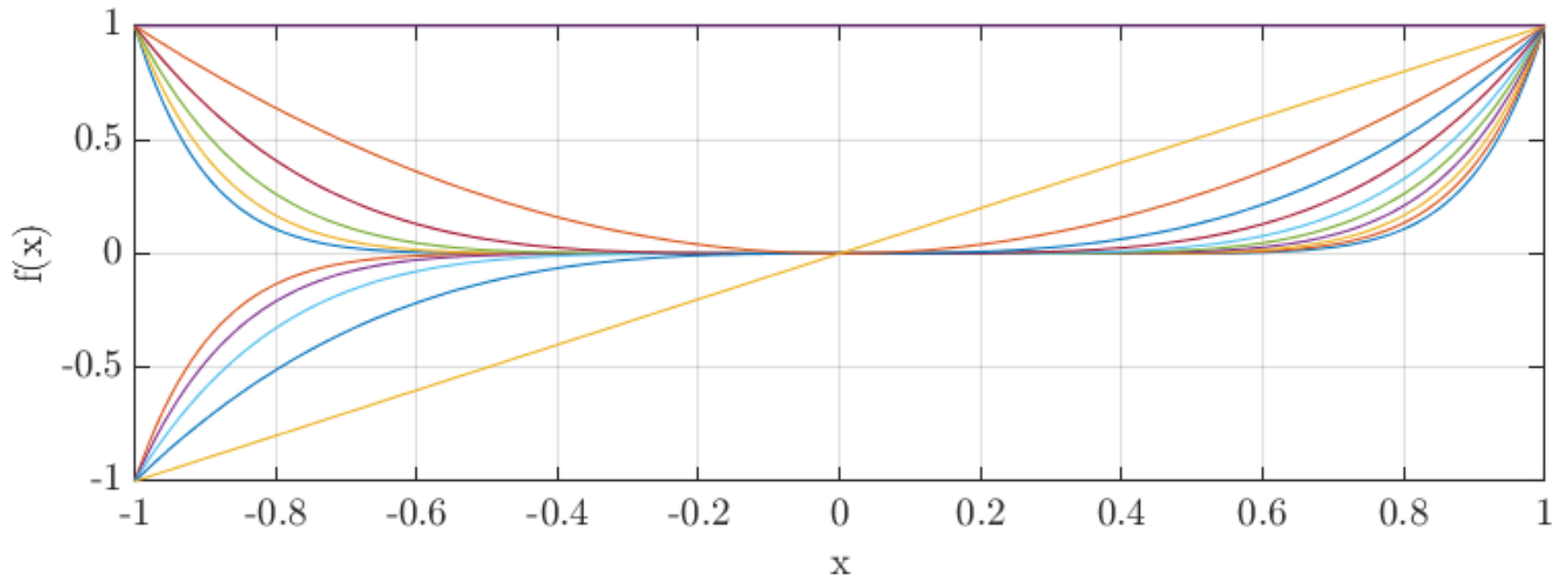
$$\mathbf{y}_m = \mathbf{B}_c \boldsymbol{\alpha}, \quad (20)$$

If $\boldsymbol{\Lambda}_\alpha$ is the covariance of $\boldsymbol{\alpha}$ (Variance of the latent space as approximation), then,

$$\boldsymbol{\Lambda}_{\mathbf{y}_m} = \mathbf{B}_c \boldsymbol{\Lambda}_\alpha \mathbf{B}_c^T. \quad (21)$$

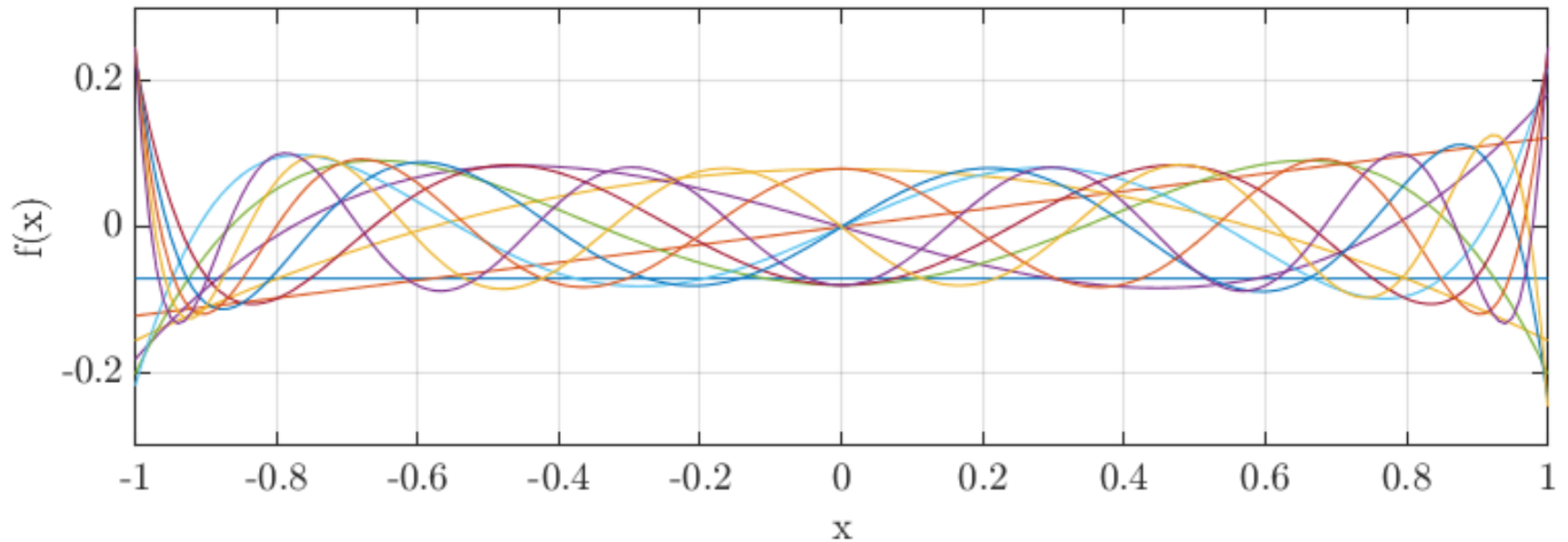


Polynomials

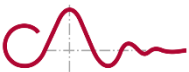
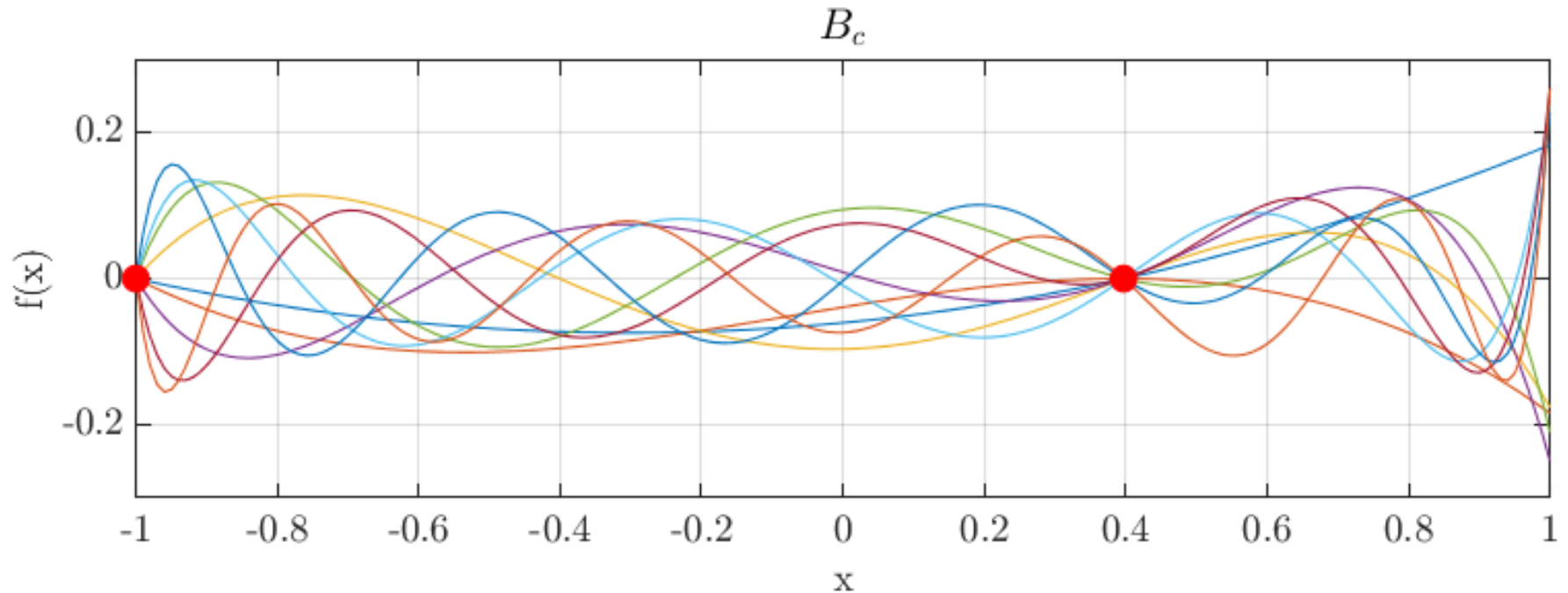


Discrete Orthogonal Polynomials

DOP



Constrained Discrete Orthogonal Polynomials



Example 2

