

# Comparing different state-of-the-art solutions for image prediction using time-series analysis

Sören Dittrich

University of Hildesheim

Summerterm 2020

# Table of contents

## 1 Introduction

## 2 Machine Learning Theory

- Image Prediction
- Autoencoder
- CNN
- RNN
- LSTM
- ConvLSTM
- Backpropagation
- BPTT

## 3 Image Prediction Architectures

- LSTM Autoencoder
- ConvLSTM Autoencoder
- Spatio-temporal Video Autoencoder
- PredNet
- PredRNN

# Contents

## 1 Introduction

## 2 Machine Learning Theory

- Image Prediction
- Autoencoder
- CNN
- RNN
- LSTM
- ConvLSTM
- Backpropagation
- BPTT

## 3 Image Prediction Architectures

- LSTM Autoencoder
- ConvLSTM Autoencoder
- Spatio-temporal Video Autoencoder
- PredNet
- PredRNN

# Contents

## 1 Introduction

## 2 Machine Learning Theory

- Image Prediction
- Autoencoder
- CNN
- RNN
- LSTM
- ConvLSTM
- Backpropagation
- BPTT

## 3 Image Prediction Architectures

- LSTM Autoencoder
- ConvLSTM Autoencoder
- Spatio-temporal Video Autoencoder
- PredNet
- PredRNN

# Image Prediction

- Field inside machine learning / computer vision

# Image Prediction

- Field inside machine learning / computer vision
- Predict future image/s, given sequence of images

# Image Prediction

- Field inside machine learning / computer vision
- Predict future image/s, given sequence of images
- $X$  the image sequence of length  $n$

# Image Prediction

- Field inside machine learning / computer vision
- Predict future image/s, given sequence of images
- $X$  the image sequence of length  $n$
- with  $X = (x_0, \dots, x_{n-1})$



# Image Prediction

- Field inside machine learning / computer vision
- Predict future image/s, given sequence of images
- $X$  the image sequence of length  $n$
- with  $X = (x_0, \dots, x_{n-1})$
- Two possible use-cases

# Image Prediction

- Field inside machine learning / computer vision
- Predict future image/s, given sequence of images
- $X$  the image sequence of length  $n$
- with  $X = (x_0, \dots, x_{n-1})$
- Two possible use-cases
  - One-frame prediction

# Image Prediction

- Field inside machine learning / computer vision
- Predict future image/s, given sequence of images
- $X$  the image sequence of length  $n$
- with  $X = (x_0, \dots, x_{n-1})$
- Two possible use-cases
  - One-frame prediction
    - Predicting  $x_n$

# Image Prediction

- Field inside machine learning / computer vision
- Predict future image/s, given sequence of images
- $X$  the image sequence of length  $n$
- with  $X = (x_0, \dots, x_{n-1})$
- Two possible use-cases
  - One-frame prediction
    - Predicting  $x_n$
  - Multi-frame prediction

# Image Prediction

- Field inside machine learning / computer vision
- Predict future image/s, given sequence of images
- $X$  the image sequence of length  $n$
- with  $X = (x_0, \dots, x_{n-1})$
- Two possible use-cases
  - One-frame prediction
    - Predicting  $x_n$
  - Multi-frame prediction
    - Predict  $t > 1$  frames into the future  $(x_n, \dots, x_{n+t-1})$

# Image Prediction

- Field inside machine learning / computer vision
- Predict future image/s, given sequence of images
- $X$  the image sequence of length  $n$
- with  $X = (x_0, \dots, x_{n-1})$
- Two possible use-cases
  - One-frame prediction
    - Predicting  $x_n$
  - Multi-frame prediction
    - Predict  $t > 1$  frames into the future  $(x_n, \dots, x_{n+t-1})$
    - Often it is one-frame prediction in a feedback loop

# Image Prediction

- Field inside machine learning / computer vision
- Predict future image/s, given sequence of images
- $X$  the image sequence of length  $n$
- with  $X = (x_0, \dots, x_{n-1})$
- Two possible use-cases
  - One-frame prediction
    - Predicting  $x_n$
  - Multi-frame prediction
    - Predict  $t > 1$  frames into the future  $(x_n, \dots, x_{n+t-1})$
    - Often it is one-frame prediction in a feedback loop
    - Propagate the error  $\rightarrow$  Greater error in later images

# Autoencoder

- Two networks chained together

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.



# Autoencoder

- Two networks chained together
  - Encoder

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$
    - $E(x) = h$

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$
    - $E(x) = h$
  - Decoder

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$
    - $E(x) = h$
  - Decoder
    - Input is  $h$

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$
    - $E(x) = h$
  - Decoder
    - Input is  $h$
    - Output is  $x'$

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$
    - $E(x) = h$
  - Decoder
    - Input is  $h$
    - Output is  $x'$
    - $D(h) = x'$

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.



# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$
    - $E(x) = h$
  - Decoder
    - Input is  $h$
    - Output is  $x'$
    - $D(h) = x'$
- Used for reconstruction  $x \approx x'$

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$
    - $E(x) = h$
  - Decoder
    - Input is  $h$
    - Output is  $x'$
    - $D(h) = x'$
- Used for reconstruction  $x \approx x'$
- Important is to prevent the network to simply copy  $x$  to  $x'$  (Interpolation)

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$
    - $E(x) = h$
  - Decoder
    - Input is  $h$
    - Output is  $x'$
    - $D(h) = x'$
- Used for reconstruction  $x \approx x'$
- Important is to prevent the network to simply copy  $x$  to  $x'$  (Interpolation)
- Simplest architecture is the undercomplete autoencoder

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$
    - $E(x) = h$
  - Decoder
    - Input is  $h$
    - Output is  $x'$
    - $D(h) = x'$
- Used for reconstruction  $x \approx x'$
- Important is to prevent the network to simply copy  $x$  to  $x'$  (Interpolation)
- Simplest architecture is the undercomplete autoencoder
  - Code smaller than input

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

- Two networks chained together
  - Encoder
    - Input is  $x$
    - Output is  $h^1$
    - $E(x) = h$
  - Decoder
    - Input is  $h$
    - Output is  $x'$
    - $D(h) = x'$
- Used for reconstruction  $x \approx x'$
- Important is to prevent the network to simply copy  $x$  to  $x'$  (Interpolation)
- Simplest architecture is the undercomplete autoencoder
  - Code smaller than input
  - Network needs to distinguish between useful and obsolete

---

<sup>1</sup> $h$  is the so named **code**. Output layer is named **bottleneck layer**.

# Autoencoder

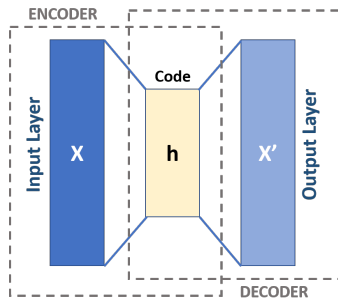


Figure: Autoencoder schema [? ]

- Convolutional Neural Network

- Convolutional Neural Network
- Consists of three stages



- Convolutional Neural Network
- Consists of three stages
  - ① Convolutional layer

- Convolutional Neural Network
- Consists of three stages
  - 1 Convolutional layer
  - 2 Non-linearity (ReLU, sigmoid, ...)

- Convolutional Neural Network
- Consists of three stages
  - 1 Convolutional layer
  - 2 Non-linearity (ReLU, sigmoid, ...)
  - 3 Pooling layer

# CNN (First stage)

- Convolutional operation is discrete

# CNN (First stage)

- Convolutional operation is discrete
- $(I * K)(i, j) = \sum_m \sum_n I(m, n) K(i - m, j - n)$

# CNN (First stage)

- Convolutional operation is discrete
- $(I * K)(i, j) = \sum_m \sum_n I(m, n)K(i - m, j - n)$

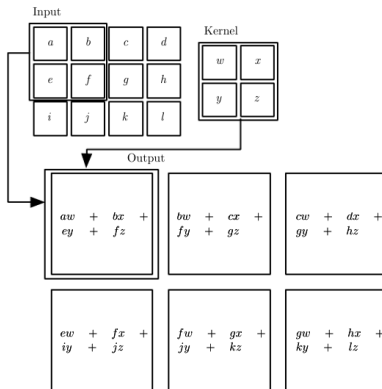


Figure: Two dimensional convolutional operation [? ]









# Backpropagation



# Contents

- 1 Introduction
- 2 Machine Learning Theory
  - Image Prediction
  - Autoencoder
  - CNN
  - RNN
  - LSTM
  - ConvLSTM
  - Backpropagation
  - BPTT
- 3 Image Prediction Architectures

- LSTM Autoencoder
- ConvLSTM Autoencoder
- Spatio-temporal Video Autoencoder
- PredNet
- PredRNN

- „Unsupervised Learning of Video Representations using LSTMs“ by Srivastava et. al. [? ]

# LSTM Autoencoder

- „Unsupervised Learning of Video Representations using LSTMs“ by Srivastava et. al. [? ]
- Using the standard LSTM from Hochreiter & Schmidhuber [? ]

# LSTM Autoencoder

- „Unsupervised Learning of Video Representations using LSTMs“ by Srivastava et. al. [? ]
- Using the standard LSTM from Hochreiter & Schmidhuber [? ]
- Autoencoder architecture

# LSTM Autoencoder

- „Unsupervised Learning of Video Representations using LSTMs“ by Srivastava et. al. [? ]
- Using the standard LSTM from Hochreiter & Schmidhuber [? ]
- Autoencoder architecture
- Useful for future image prediction & image reconstruction



# LSTM Autoencoder

- „Unsupervised Learning of Video Representations using LSTMs“ by Srivastava et. al. [?] ]
- Using the standard LSTM from Hochreiter & Schmidhuber [?] ]
- Autoencoder architecture
- Useful for future image prediction & image reconstruction
- Typical baseline for newer, more advanced algorithms

# LSTM Autoencoder

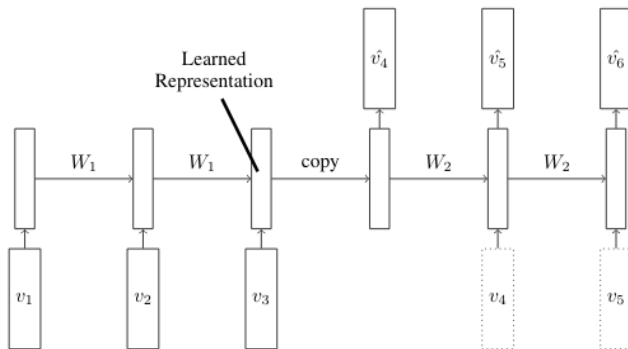


Figure: Future image prediction model [? ]

# LSTM Autoencoder

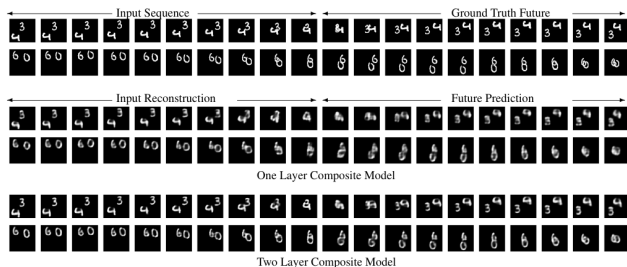


Figure: Results of MovingMNIST experiment [? ]

- „Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting“ by Shi et. al. [?] ]

- „Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting“ by Shi et. al. [? ]
- Similar to LSTM Autoencoder, but uses ConvLSTM instead

# ConvLSTM Autoencoder

- „Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting“ by Shi et. al. [? ]
- Similar to LSTM Autoencoder, but uses ConvLSTM instead
- Outperforms the LSTM Autoencoder

# ConvLSTM Autoencoder

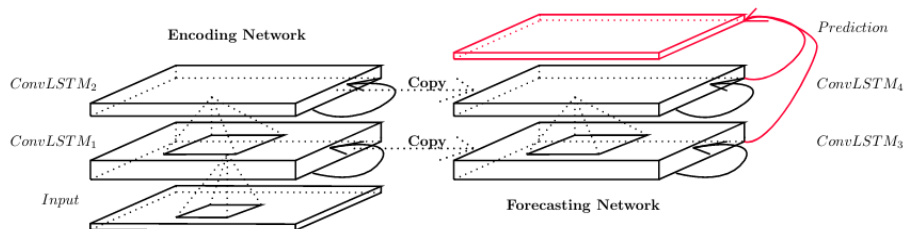


Figure: Future image prediction model [? ]

# ConvLSTM Autoencoder

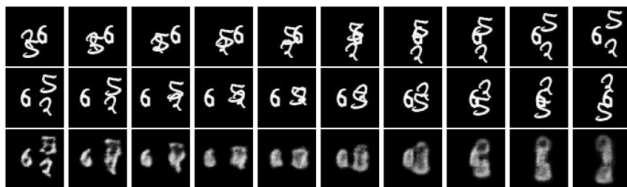


Figure: Results of MovingMNIST experiment [? ]



# Spatio-temporal Video Autoencoder

- „Spatio-Temporal Video Autoencoder With Differentiable Memory “by Patraucean et. al. [? ]

# Spatio-temporal Video Autoencoder

- „Spatio-Temporal Video Autoencoder With Differentiable Memory “by Patraucean et. al. [? ]
-

# Spatio-temporal Video Autoencoder

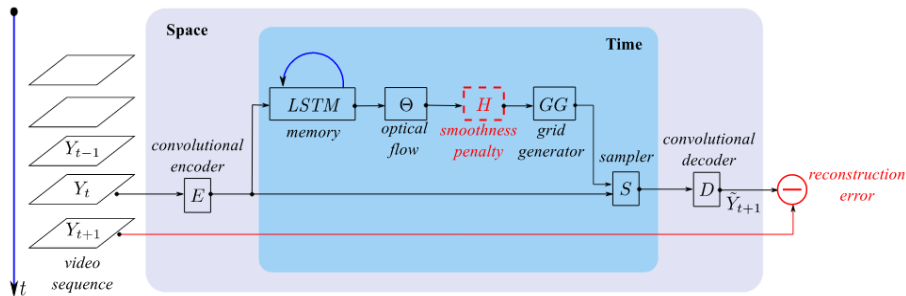


Figure: Spatio-temporal Video Autoencoder Architecture [? ]

# Spatio-temporal Video Autoencoder



Figure: Results of MovingMNIST experiment [? ]



