

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 Scientific questions
- 2 Introduction
- 3 Machine Learning Theory
 - Image Prediction
 - Autoencoder
 - CNN
 - RNN
 - LSTM
 - ConvLSTM
 - Backpropagation
 - BPTT
- 4 Image Prediction Architectures
 - LSTM Autoencoder
 - ConvLSTM Autoencoder
 - Spatio-temporal Video Autoencoder
 - PredNet
 - PredRNN
- 5 Experiments
 - Experimental setup
 - Experimental results (First setup)

Contents

- 1 Scientific questions
- 2 Introduction
- 3 Machine Learning Theory
 - Image Prediction
 - Autoencoder
 - CNN
 - RNN
 - LSTM
 - ConvLSTM
 - Backpropagation
 - BPTT
- 4 Image Prediction Architectures
 - LSTM Autoencoder
 - ConvLSTM Autoencoder
 - Spatio-temporal Video Autoencoder
 - PredNet
 - PredRNN
- 5 Experiments
 - Experimental setup
 - Experimental results (First setup)

- 1 What different types of image-/video-prediction architectures exist?

- ① What different types of image-/video-prediction architectures exist?
- ② How important is the choice of the recurrent module for the runtime and performance of the algorithm?

- 1 What different types of image-/video-prediction architectures exist?
- 2 How important is the choice of the recurrent module for the runtime and performance of the algorithm?
- 3 How important is hyperparameter optimization?

Contents

- 1 Scientific questions
- 2 Introduction
- 3 Machine Learning Theory
 - Image Prediction
 - Autoencoder
 - CNN
 - RNN
 - LSTM
 - ConvLSTM
 - Backpropagation
 - BPTT
- 4 Image Prediction Architectures
 - LSTM Autoencoder
 - ConvLSTM Autoencoder
 - Spatio-temporal Video Autoencoder
 - PredNet
 - PredRNN
- 5 Experiments
 - Experimental setup
 - Experimental results (First setup)

- Overview of necessary machine learning theory

- Overview of necessary machine learning theory
- Overview of different state-of-the-art solutions for image-/video-prediction

Introduction

- Overview of necessary machine learning theory
- Overview of different state-of-the-art solutions for image-/video-prediction
- Implementation of three image-/video-prediction algorithms in PyTorch [10]

- Overview of necessary machine learning theory
- Overview of different state-of-the-art solutions for image-/video-prediction
- Implementation of three image-/video-prediction algorithms in PyTorch [10]
- Change of recurrent module during experiments to achieve performance improvement

- Overview of necessary machine learning theory
- Overview of different state-of-the-art solutions for image-/video-prediction
- Implementation of three image-/video-prediction algorithms in PyTorch [10]
- Change of recurrent module during experiments to achieve performance improvement
- Discussion of the overview and the performed experiments

- Overview of necessary machine learning theory
- Overview of different state-of-the-art solutions for image-/video-prediction
- Implementation of three image-/video-prediction algorithms in PyTorch [10]
- Change of recurrent module during experiments to achieve performance improvement
- Discussion of the overview and the performed experiments
- Close with a conclusion

Contents

- 1 Scientific questions
- 2 Introduction
- 3 Machine Learning Theory
 - Image Prediction
 - Autoencoder
 - CNN
 - RNN
 - LSTM
 - ConvLSTM
 - Backpropagation
 - BPTT
- 4 Image Prediction Architectures
 - LSTM Autoencoder
 - ConvLSTM Autoencoder
 - Spatio-temporal Video Autoencoder
 - PredNet
 - PredRNN
- 5 Experiments
 - Experimental setup
 - Experimental results (First setup)

- 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

- Two networks chained together

¹ h is the so named **code**. Output layer is named **bottleneck layer**.

Autoencoder

- Two networks chained together
 - Encoder

¹ h is the so named **code**. Output layer is named **bottleneck layer**.

Autoencoder

- Two networks chained together
 - Encoder
 - Input is x

¹ 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

¹ 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
 - $E(x) = h$

¹ 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

¹ 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

¹ h is the so named **code**. Output layer is named **bottleneck layer**.

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

¹ h is the so named **code**. Output layer is named **bottleneck layer**.

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

¹ h is the so named **code**. Output layer is named **bottleneck layer**.

- 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 $x \approx x'$

¹ h is the so named **code**. Output layer is named **bottleneck layer**.

- 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)

¹ h is the so named **code**. Output layer is named **bottleneck layer**.

- 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

¹ h is the so named **code**. Output layer is named **bottleneck layer**.

- 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

¹ h is the so named **code**. Output layer is named **bottleneck layer**.

- 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

¹ h is the so named **code**. Output layer is named **bottleneck layer**.

Autoencoder

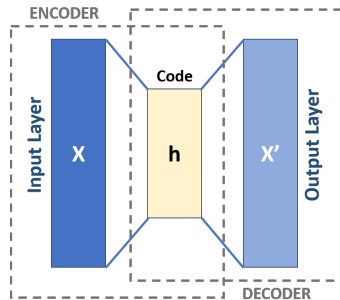


Figure: Autoencoder schema [9]

- 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

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

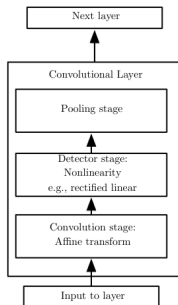


Figure: Stages of a CNN [4]

CNN (First stage)

- Discrete convolutional operation

CNN (First stage)

- Discrete convolutional operation
- $(I * K)(i, j) = \sum_m \sum_n I(m, n) K(i - m, j - n)$

CNN (First stage)

- Discrete convolutional operation
- $(I * K)(i, j) = \sum_m \sum_n I(m, n)K(i - m, j - n)$

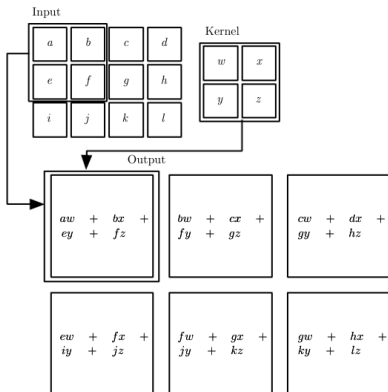


Figure: Two dimensional convolutional operation [4]

CNN (Second stage)

- Non-linearity layer

CNN (Second stage)

- Non-linearity layer
- ReLU: $R(x) = \max(0, x)$

CNN (Second stage)

- Non-linearity layer
- ReLU: $R(x) = \max(0, x)$
- Sigmoid: $\sigma(x) = \frac{1}{1+e^{-x}}$

CNN (Third stage)

- Pooling layer

CNN (Third stage)

- Pooling layer
- Aggregate output / Making output invariant to input

CNN (Third stage)

- Pooling layer
- Aggregate output / Making output invariant to input
- E.g. Max-Pooling

CNN (Third stage)

- Pooling layer
- Aggregate output / Making output invariant to input
- E.g. Max-Pooling

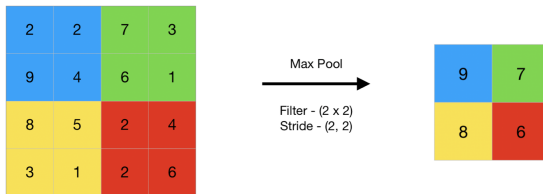


Figure: Max-Pooling (Filter := Kernel)[1]

- Recurrent Neural Network

- Recurrent Neural Network
- Handling sequential data

- Recurrent Neural Network
- Handling sequential data
- $X = (x_0, \dots, x_{t-1}), |X| = t$

- Recurrent Neural Network
- Handling sequential data
- $X = (x_0, \dots, x_{t-1}), |X| = t$
- Used for e.g. time-series analysis, image-/video-prediction

- Recurrent Neural Network
- Handling sequential data
- $X = (x_0, \dots, x_{t-1})$, $|X| = t$
- Used for e.g. time-series analysis, image-/video-prediction
- Gets output of last iteration as „additional“input

- Recurrent Neural Network
- Handling sequential data
- $X = (x_0, \dots, x_{t-1})$, $|X| = t$
- Used for e.g. time-series analysis, image-/video-prediction
- Gets output of last iteration as „additional“input
- $\hat{y}^t = f_{\theta}(\hat{y}^{t-1}; x^t) = f_{\theta}(f_{\theta}(\hat{y}^{t-2}; x^{t-1}); x^t) = \dots$

- Recurrent Neural Network
- Handling sequential data
- $X = (x_0, \dots, x_{t-1})$, $|X| = t$
- Used for e.g. time-series analysis, image-/video-prediction
- Gets output of last iteration as „additional“input
- $\hat{y}^t = f_{\theta}(\hat{y}^{t-1}; x^t) = f_{\theta}(f_{\theta}(\hat{y}^{t-2}; x^{t-1}); x^t) = \dots$
- Requires \hat{y}^0 to be initialized (Often done with zero)

- Recurrent Neural Network
- Handling sequential data
- $X = (x_0, \dots, x_{t-1})$, $|X| = t$
- Used for e.g. time-series analysis, image-/video-prediction
- Gets output of last iteration as „additional“input
- $\hat{y}^t = f_{\theta}(\hat{y}^{t-1}; x^t) = f_{\theta}(f_{\theta}(\hat{y}^{t-2}; x^{t-1}); x^t) = \dots$
- Requires \hat{y}^0 to be initialized (Often done with zero)
- Typically trained using BPTT

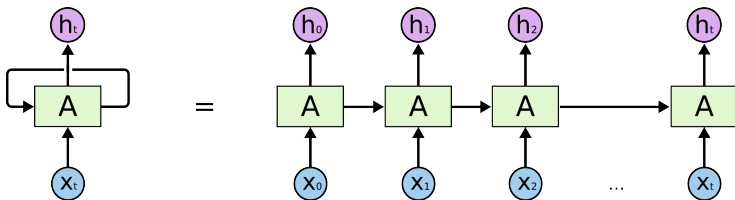


Figure: RNN schema. **Left:** Folded graph, **Right:** Unfolded graph [2]

- Long Short-term Memory

LSTM

- Long Short-term Memory
- Invented by Hochreiter and Schmidhuber [5]

- Long Short-term Memory
- Invented by Hochreiter and Schmidhuber [5]
- Avoids critical problem of standard RNN: Saving long-term dependencies

- Long Short-term Memory
- Invented by Hochreiter and Schmidhuber [5]
- Avoids critical problem of standard RNN: Saving long-term dependencies

$$i_t = \sigma(w_{x_i} x_t + w_{h_i} h_{t-1} + b_i) \quad (1)$$

- Long Short-term Memory
- Invented by Hochreiter and Schmidhuber [5]
- Avoids critical problem of standard RNN: Saving long-term dependencies

$$i_t = \sigma(w_{x_i}x_t + w_{h_i}h_{t-1} + b_i) \quad (1)$$

$$f_t = \sigma(w_{x_f}x_t + w_{h_f}h_{t-1} + b_f) \quad (2)$$

- Long Short-term Memory
- Invented by Hochreiter and Schmidhuber [5]
- Avoids critical problem of standard RNN: Saving long-term dependencies

$$i_t = \sigma(w_{x_i}x_t + w_{h_i}h_{t-1} + b_i) \quad (1)$$

$$f_t = \sigma(w_{x_f}x_t + w_{h_f}h_{t-1} + b_f) \quad (2)$$

$$c_t = f_t c_{t-1} + i_t \tanh(w_{x_c}x_t + w_{h_c}h_{t-1} + b_c) \quad (3)$$

- Long Short-term Memory
- Invented by Hochreiter and Schmidhuber [5]
- Avoids critical problem of standard RNN: Saving long-term dependencies

$$i_t = \sigma(w_{x_i}x_t + w_{h_i}h_{t-1} + b_i) \quad (1)$$

$$f_t = \sigma(w_{x_f}x_t + w_{h_f}h_{t-1} + b_f) \quad (2)$$

$$c_t = f_t c_{t-1} + i_t \tanh(w_{x_c}x_t + w_{h_c}h_{t-1} + b_c) \quad (3)$$

$$o_t = \sigma(w_{x_o}x_t + w_{h_o}h_{t-1} + b_o) \quad (4)$$

- Long Short-term Memory
- Invented by Hochreiter and Schmidhuber [5]
- Avoids critical problem of standard RNN: Saving long-term dependencies

$$i_t = \sigma(w_{x_i}x_t + w_{h_i}h_{t-1} + b_i) \quad (1)$$

$$f_t = \sigma(w_{x_f}x_t + w_{h_f}h_{t-1} + b_f) \quad (2)$$

$$c_t = f_t c_{t-1} + i_t \tanh(w_{x_c}x_t + w_{h_c}h_{t-1} + b_c) \quad (3)$$

$$o_t = \sigma(w_{x_o}x_t + w_{h_o}h_{t-1} + b_o) \quad (4)$$

$$h_t = o_t \tanh(c_t) \quad (5)$$

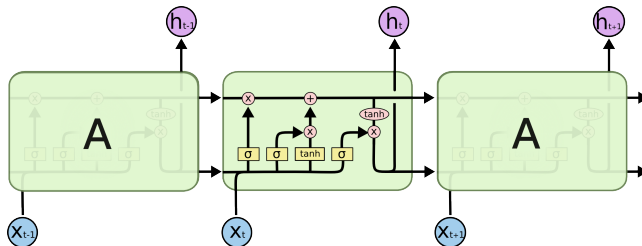


Figure: LSTM Architecture [2]

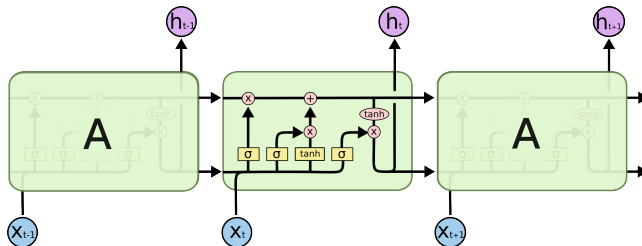


Figure: LSTM Architecture [2]

- Top horizontal line: Cell-state

- Invented by Shi et al. [14]

ConvLSTM

- Invented by Shi et al. [14]
- LSTM with peephole using convlayer

- Invented by Shi et al. [14]
- LSTM with peephole using convlayer

$$i_t = \sigma(x_t * w_{x_i} + h_{t-1} * w_{h_i} + w_{i_b}) \quad (6)$$

- Invented by Shi et al. [14]
- LSTM with peephole using convlayer

$$i_t = \sigma(x_t * w_{x_i} + h_{t-1} * w_{h_i} + w_{i_b}) \quad (6)$$

$$f_t = \sigma(x_t * w_{x_f} + h_{t-1} * w_{h_f} + w_{f_b}) \quad (7)$$

- Invented by Shi et al. [14]
- LSTM with peephole using convlayer

$$i_t = \sigma(x_t * w_{x_i} + h_{t-1} * w_{h_i} + w_{i_b}) \quad (6)$$

$$f_t = \sigma(x_t * w_{x_f} + h_{t-1} * w_{h_f} + w_{f_b}) \quad (7)$$

$$\tilde{c}_t = \tanh(x_t * w_{x_{\tilde{c}}} + h_{t-1} * w_{h_{\tilde{c}}} + w_{\tilde{c}_b}) \quad (8)$$

- Invented by Shi et al. [14]
- LSTM with peephole using convlayer

$$i_t = \sigma(x_t * w_{x_i} + h_{t-1} * w_{h_i} + w_{i_b}) \quad (6)$$

$$f_t = \sigma(x_t * w_{x_f} + h_{t-1} * w_{h_f} + w_{f_b}) \quad (7)$$

$$\tilde{c}_t = \tanh(x_t * w_{x_{\tilde{c}}} + h_{t-1} * w_{h_{\tilde{c}}} + w_{\tilde{c}_b}) \quad (8)$$

$$c_t = \tilde{c}_t \odot i_t + c_{t-1} \odot f_t \quad (9)$$

- Invented by Shi et al. [14]
- LSTM with peephole using convlayer

$$i_t = \sigma(x_t * w_{x_i} + h_{t-1} * w_{h_i} + w_{i_b}) \quad (6)$$

$$f_t = \sigma(x_t * w_{x_f} + h_{t-1} * w_{h_f} + w_{f_b}) \quad (7)$$

$$\tilde{c}_t = \tanh(x_t * w_{x_{\tilde{c}}} + h_{t-1} * w_{h_{\tilde{c}}} + w_{\tilde{c}_b}) \quad (8)$$

$$c_t = \tilde{c}_t \odot i_t + c_{t-1} \odot f_t \quad (9)$$

$$o_t = \sigma(x_t * w_{x_o} + h_{t-1} * w_{h_o} + w_{o_b}) \quad (10)$$

- Invented by Shi et al. [14]
- LSTM with peephole using convlayer

$$i_t = \sigma(x_t * w_{x_i} + h_{t-1} * w_{h_i} + w_{i_b}) \quad (6)$$

$$f_t = \sigma(x_t * w_{x_f} + h_{t-1} * w_{h_f} + w_{f_b}) \quad (7)$$

$$\tilde{c}_t = \tanh(x_t * w_{x_{\tilde{c}}} + h_{t-1} * w_{h_{\tilde{c}}} + w_{\tilde{c}_b}) \quad (8)$$

$$c_t = \tilde{c}_t \odot i_t + c_{t-1} \odot f_t \quad (9)$$

$$o_t = \sigma(x_t * w_{x_o} + h_{t-1} * w_{h_o} + w_{o_b}) \quad (10)$$

$$h_t = o_t \odot \tanh(c_t) \quad (11)$$

Backpropagation

- Backpropagation was invented by Rumelhart et al. [13] in 1986.

Backpropagation

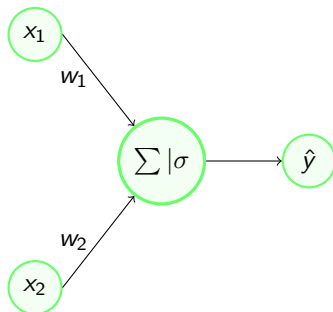
- Backpropagation was invented by Rumelhart et al. [13] in 1986.
- Most used algorithm for training neural nets.

Backpropagation

- Backpropagation was invented by Rumelhart et al. [13] in 1986.
- Most used algorithm for training neural nets.
- Simple example on Perceptron [12].

Backpropagation

- Backpropagation was invented by Rumelhart et al. [13] in 1986.
- Most used algorithm for training neural nets.
- Simple example on Perceptron [12].



Backpropagation Forward Pass

Forward Pass:

Backpropagation Forward Pass

Forward Pass:



$$\hat{y} = \sigma\left(\sum_{i=1}^2 x_i w_i\right) \quad (12)$$

Backpropagation Forward Pass

Forward Pass:

-

$$\hat{y} = \sigma\left(\sum_{i=1}^2 x_i w_i\right) \quad (12)$$

- σ is the sigmoid function

Backpropagation Forward Pass

Forward Pass:

-

$$\hat{y} = \sigma\left(\sum_{i=1}^2 x_i w_i\right) \quad (12)$$

- σ is the sigmoid function

-

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (13)$$

Backpropagation Forward Pass

Forward Pass:



$$\hat{y} = \sigma\left(\sum_{i=1}^2 x_i w_i\right) \quad (12)$$

- σ is the sigmoid function



$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (13)$$

- Computing the error with e.g. **MSE**

Backpropagation Forward Pass

Forward Pass:

-

$$\hat{y} = \sigma\left(\sum_{i=1}^2 x_i w_i\right) \quad (12)$$

- σ is the sigmoid function

-

$$\sigma(x) = \frac{1}{1 + e^{-x}} \quad (13)$$

- Computing the error with e.g. **MSE**

-

$$L(y, \hat{y}) = \frac{1}{N} \|y - \hat{y}\|_2^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (14)$$

Backpropagation Backward Pass

Backward Pass:

Backpropagation Backward Pass

Backward Pass:



$$\frac{\partial L}{\partial \hat{y}} = \frac{2}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \quad (15)$$

Backpropagation Backward Pass

Backward Pass:

-

$$\frac{\partial L}{\partial \hat{y}} = \frac{2}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \quad (15)$$

-

$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x)) \quad (16)$$

Backpropagation Backward Pass

Backward Pass:

- $$\frac{\partial L}{\partial \hat{y}} = \frac{2}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \quad (15)$$

- $$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x)) \quad (16)$$

- $$\frac{\partial L}{\partial \sum_{i=1}^2 x_i w_i} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sum_{i=1}^2 x_i w_i} = \frac{2}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \cdot \sigma\left(\sum_{i=1}^2 x_i w_i\right) (1 - \sigma\left(\sum_{i=1}^2 x_i w_i\right)) \quad (17)$$

Backpropagation Backward Pass

Backward Pass:

- $$\frac{\partial L}{\partial \hat{y}} = \frac{2}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \quad (15)$$

- $$\frac{\partial \sigma(x)}{\partial x} = \sigma(x)(1 - \sigma(x)) \quad (16)$$

- $$\frac{\partial L}{\partial \sum_{i=1}^2 x_i w_i} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial \sum_{i=1}^2 x_i w_i} = \frac{2}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \cdot \sigma\left(\sum_{i=1}^2 x_i w_i\right) (1 - \sigma\left(\sum_{i=1}^2 x_i w_i\right)) \quad (17)$$

- $$\frac{\partial L}{\partial w_1} = \dots = \frac{2}{N} \sum_{i=1}^N (y_i - \hat{y}_i) \cdot \sigma\left(\sum_{i=1}^2 x_i w_i\right) (1 - \sigma\left(\sum_{i=1}^2 x_i w_i\right)) \cdot x_1 \quad (18)$$

Backpropagation Update

Update weights using Gradient Descent

Backpropagation Update

Update weights using Gradient Descent



$$w_1 = w_1^{old} - \lambda \cdot \frac{\partial L}{\partial w_1} \quad (19)$$

Backpropagation Update

Update weights using Gradient Descent



$$w_1 = w_1^{old} - \lambda \cdot \frac{\partial L}{\partial w_1} \quad (19)$$

- λ is the learning rate

Backpropagation Update

Update weights using Gradient Descent

-

$$w_1 = w_1^{old} - \lambda \cdot \frac{\partial L}{\partial w_1} \quad (19)$$

- λ is the learning rate
- All steps are performed iterative, until convergence

Backpropagation Update

Update weights using Gradient Descent

- $$w_1 = w_1^{old} - \lambda \cdot \frac{\partial L}{\partial w_1} \quad (19)$$

- λ is the learning rate
- All steps are performed iterative, until convergence
- It would be also possible to use e.g. Newton instead of Gradient Descent

Update weights using Gradient Descent

-

$$w_1 = w_1^{old} - \lambda \cdot \frac{\partial L}{\partial w_1} \quad (19)$$

- λ is the learning rate
- All steps are performed iterative, until convergence
- It would be also possible to use e.g. Newton instead of Gradient Descent
 - Gradient Descent is used more, because of its simplicity

Backpropagation Update

Update weights using Gradient Descent

-

$$w_1 = w_1^{old} - \lambda \cdot \frac{\partial L}{\partial w_1} \quad (19)$$

- λ is the learning rate
- All steps are performed iterative, until convergence
- It would be also possible to use e.g. Newton instead of Gradient Descent
 - Gradient Descent is used more, because of its simplicity
 - and because of its parallelism properties [13]

- Backpropagation through time

- Backpropagation through time
- Invented by Paul Werbos [17]

- Backpropagation through time
- Invented by Paul Werbos [17]
- Unfolding graph during backpass and propagate through the steps

- Backpropagation through time
- Invented by Paul Werbos [17]
- Unfolding graph during backpass and propagate through the steps
- Simple example on RNN found in Goodfellow [4]

Forward Pass:

Forward Pass:



$$h_t = \tanh(Wh_{t-1} + Ux_t + b_1) \quad (20)$$

$$o_t = Vh_t + b_2 \quad (21)$$

$$\hat{y}_t = \sigma(o_t) \quad (22)$$

Forward Pass:



$$h_t = \tanh(Wh_{t-1} + Ux_t + b_1) \quad (20)$$

$$o_t = Vh_t + b_2 \quad (21)$$

$$\hat{y}_t = \sigma(o_t) \quad (22)$$

- U, V, W weight matrices for input-to-hidden, hidden-to-output, hidden-to-hidden

Forward Pass:



$$h_t = \tanh(Wh_{t-1} + Ux_t + b_1) \quad (20)$$

$$o_t = Vh_t + b_2 \quad (21)$$

$$\hat{y}_t = \sigma(o_t) \quad (22)$$

- U, V, W weight matrices for input-to-hidden, hidden-to-output, hidden-to-hidden
- Computing the error with e.g. **MSE** for every iteration

Forward Pass:



$$h_t = \tanh(Wh_{t-1} + Ux_t + b_1) \quad (20)$$

$$o_t = Vh_t + b_2 \quad (21)$$

$$\hat{y}_t = \sigma(o_t) \quad (22)$$

- U, V, W weight matrices for input-to-hidden, hidden-to-output, hidden-to-hidden
- Computing the error with e.g. **MSE** for every iteration



$$L(y, \hat{y}) = \sum_t (L_t(y_t, \hat{y}_t)) \quad (23)$$

BPTT Backward Pass

Backward Pass:

BPTT Backward Pass

Backward Pass:



$$\frac{\partial L}{\partial V} = \sum_t \frac{\partial L_t}{\partial V} \quad (24)$$

BPTT Backward Pass

Backward Pass:

- $$\frac{\partial L}{\partial V} = \sum_t \frac{\partial L_t}{\partial V} \quad (24)$$

- $$\frac{\partial L_t}{\partial \hat{y}_t} = \frac{2}{N} \sum_{i=1}^N (y_t^i - \hat{y}_t^i) \quad (25)$$

Backward Pass:

- $$\frac{\partial L}{\partial V} = \sum_t \frac{\partial L_t}{\partial V} \quad (24)$$

- $$\frac{\partial L_t}{\partial \hat{y}_t} = \frac{2}{N} \sum_{i=1}^N (y_t^i - \hat{y}_t^i) \quad (25)$$

- $$\frac{\partial L_t}{\partial o_t} = \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial o_t} = \frac{2}{N} \sum_{i=1}^N (y_t^i - \hat{y}_t^i) \cdot \sigma(o_t)(1 - \sigma(o_t)) \quad (26)$$

Backward Pass:

- $$\frac{\partial L}{\partial V} = \sum_t \frac{\partial L_t}{\partial V} \quad (24)$$

- $$\frac{\partial L_t}{\partial \hat{y}_t} = \frac{2}{N} \sum_{i=1}^N (y_t^i - \hat{y}_t^i) \quad (25)$$

- $$\frac{\partial L_t}{\partial o_t} = \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial o_t} = \frac{2}{N} \sum_{i=1}^N (y_t^i - \hat{y}_t^i) \cdot \sigma(o_t)(1 - \sigma(o_t)) \quad (26)$$

- $$\frac{\partial L_t}{\partial V} = \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial V} = \frac{2}{N} \sum_{i=1}^N (y_t^i - \hat{y}_t^i) \cdot \sigma(o_t)(1 - \sigma(o_t)) \cdot h_t \quad (27)$$

Backward Pass:

- $$\frac{\partial L}{\partial V} = \sum_t \frac{\partial L_t}{\partial V} \quad (24)$$

- $$\frac{\partial L_t}{\partial \hat{y}_t} = \frac{2}{N} \sum_{i=1}^N (y_t^i - \hat{y}_t^i) \quad (25)$$

- $$\frac{\partial L_t}{\partial o_t} = \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial o_t} = \frac{2}{N} \sum_{i=1}^N (y_t^i - \hat{y}_t^i) \cdot \sigma(o_t)(1 - \sigma(o_t)) \quad (26)$$

- $$\frac{\partial L_t}{\partial V} = \frac{\partial L_t}{\partial \hat{y}_t} \cdot \frac{\partial \hat{y}_t}{\partial o_t} \cdot \frac{\partial o_t}{\partial V} = \frac{2}{N} \sum_{i=1}^N (y_t^i - \hat{y}_t^i) \cdot \sigma(o_t)(1 - \sigma(o_t)) \cdot h_t \quad (27)$$

- $$\frac{\partial L}{\partial V} = \sum_t \frac{\partial L_t}{\partial V} \quad (28)$$

Update weights using Gradient Descent

Update weights using Gradient Descent

- $$V = V^{old} - \lambda \cdot \frac{\partial L}{\partial V} \quad (29)$$

Contents

- 1 Scientific questions
- 2 Introduction
- 3 Machine Learning Theory
 - Image Prediction
 - Autoencoder
 - CNN
 - RNN
 - LSTM
 - ConvLSTM
 - Backpropagation
 - BPTT
- 4 Image Prediction Architectures
 - LSTM Autoencoder
 - ConvLSTM Autoencoder
 - Spatio-temporal Video Autoencoder
 - PredNet
 - PredRNN
- 5 Experiments
 - Experimental setup
 - Experimental results (First setup)

- Description of five state-of-the-art image-/video-prediction architectures

- Description of five state-of-the-art image-/video-prediction architectures
- Different approaches, but all use a LSTM as recurrent sub-module

- Description of five state-of-the-art image-/video-prediction architectures
- Different approaches, but all use a LSTM as recurrent sub-module
- Experiments are shown on MovingMNIST [7]

- „Unsupervised Learning of Video Representations using LSTMs“ by Srivastava et. al. [15]

- „Unsupervised Learning of Video Representations using LSTMs“ by Srivastava et. al. [15]
- Using the standard LSTM from Hochreiter & Schmidhuber [5]

- „Unsupervised Learning of Video Representations using LSTMs“ by Srivastava et. al. [15]
- Using the standard LSTM from Hochreiter & Schmidhuber [5]
- Autoencoder architecture

- „Unsupervised Learning of Video Representations using LSTMs“ by Srivastava et. al. [15]
- Using the standard LSTM from Hochreiter & Schmidhuber [5]
- Autoencoder architecture
- Used for future image prediction & image reconstruction

- „Unsupervised Learning of Video Representations using LSTMs“ by Srivastava et. al. [15]
- Using the standard LSTM from Hochreiter & Schmidhuber [5]
- Autoencoder architecture
- Used for future image prediction & image reconstruction
- Typical baseline for newer, more advanced algorithms

LSTM Autoencoder

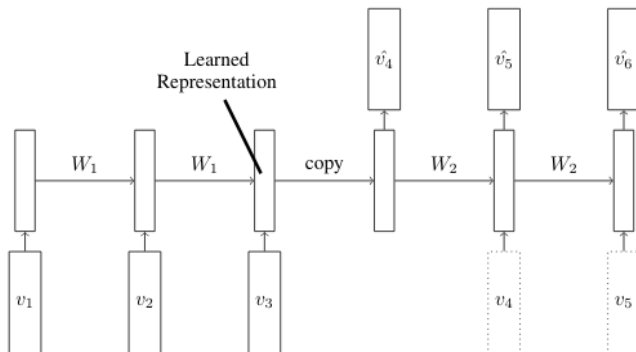


Figure: Future image prediction model [15]

- Composed model can do image prediction and reconstruction simultaneously

- Composed model can do image prediction and reconstruction simultaneously
 - Consisting of one encoder and two decoder

- Composed model can do image prediction and reconstruction simultaneously
 - Consisting of one encoder and two decoder
 - Good example of representation learning algorithms

- Composed model can do image prediction and reconstruction simultaneously
 - Consisting of one encoder and two decoder
 - Good example of representation learning algorithms
 - Both use the same encoder but differently trained decoder

- Composed model can do image prediction and reconstruction simultaneously
 - Consisting of one encoder and two decoder
 - Good example of representation learning algorithms
 - Both use the same encoder but differently trained decoder
- End-to-End differentiable

- Composed model can do image prediction and reconstruction simultaneously
 - Consisting of one encoder and two decoder
 - Good example of representation learning algorithms
 - Both use the same encoder but differently trained decoder
- End-to-End differentiable
- Trained with BPTT and CE-Loss for MovingMNIST (two digits per image)

- Composed model can do image prediction and reconstruction simultaneously
 - Consisting of one encoder and two decoder
 - Good example of representation learning algorithms
 - Both use the same encoder but differently trained decoder
- End-to-End differentiable
- Trained with BPTT and CE-Loss for MovingMNIST (two digits per image)
- Used for multi-frame prediction 10 input images \Rightarrow 10 output images

LSTM Autoencoder

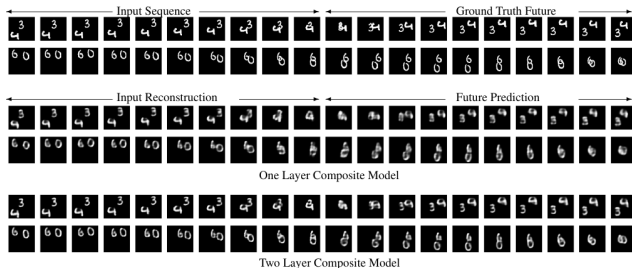


Figure: Results of MovingMNIST experiment [15]

- „Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting“ by Shi et. al. [14]

- „Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting“ by Shi et. al. [14]
- Similar to LSTM Autoencoder, but uses novel ConvLSTM instead

- „Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting“ by Shi et. al. [14]
- Similar to LSTM Autoencoder, but uses novel ConvLSTM instead
- ConvLSTM with peephole connection

- „Convolutional LSTM Network: A Machine Learning Approach for Precipitation Nowcasting“ by Shi et. al. [14]
- Similar to LSTM Autoencoder, but uses novel ConvLSTM instead
- ConvLSTM with peephole connection
- Outperforms the LSTM Autoencoder „Captures spatiotemporal correlations better“

ConvLSTM Autoencoder

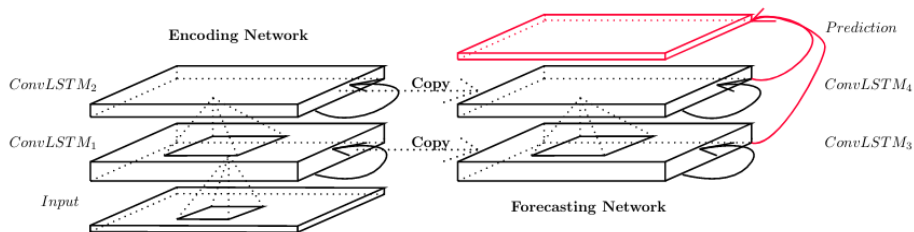


Figure: Future image prediction model [14]

- End-to-End differentiable

- End-to-End differentiable
- Trained using BPTT with CE-Loss for MovingMNIST (Three digits per image)

- End-to-End differentiable
- Trained using BPTT with CE-Loss for MovingMNIST (Three digits per image)
- Used for multi-frame prediction 10 input images \Rightarrow 10 output images

- End-to-End differentiable
- Trained using BPTT with CE-Loss for MovingMNIST (Three digits per image)
- Used for multi-frame prediction 10 input images \Rightarrow 10 output images
- Architecture is only used for forecasting, but might be able to have the same composed model as Srivastava et al. solution

ConvLSTM Autoencoder

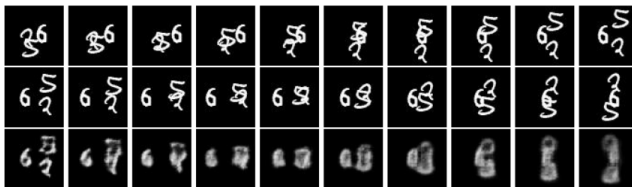


Figure: Results of MovingMNIST experiment [15]

- „Spatio-Temporal Video Autoencoder With Differentiable Memory “by Patraucean et. al. [11]

Spatio-temporal Video Autoencoder

- „Spatio-Temporal Video Autoencoder With Differentiable Memory “by Patraucean et. al. [11]
- Using ConvLSTM as recurrent sub-module without peephole connection

Spatio-temporal Video Autoencoder

- „Spatio-Temporal Video Autoencoder With Differentiable Memory “by Patraucean et. al. [11]
- Using ConvLSTM as recurrent sub-module without peephole connection
- More complex architecture, by nesting an autoencoder inside another

Spatio-temporal Video Autoencoder

- „Spatio-Temporal Video Autoencoder With Differentiable Memory “by Patraucean et. al. [11]
- Using ConvLSTM as recurrent sub-module without peephole connection
- More complex architecture, by nesting an autoencoder inside another
 - The inner autoencoder is for temporal information

- „Spatio-Temporal Video Autoencoder With Differentiable Memory “by Patraucean et. al. [11]
- Using ConvLSTM as recurrent sub-module without peephole connection
- More complex architecture, by nesting an autoencoder inside another
 - The inner autoencoder is for temporal information
 - The outer autoencoder is for spatial information

Spatio-temporal Video Autoencoder

- „Spatio-Temporal Video Autoencoder With Differentiable Memory “by Patraucean et. al. [11]
- Using ConvLSTM as recurrent sub-module without peephole connection
- More complex architecture, by nesting an autoencoder inside another
 - The inner autoencoder is for temporal information
 - The outer autoencoder is for spatial information
- Spatial autoencoder is a simple undercomplete autoencoder

Spatio-temporal Video Autoencoder

- „Spatio-Temporal Video Autoencoder With Differentiable Memory “by Patraucean et. al. [11]
- Using ConvLSTM as recurrent sub-module without peephole connection
- More complex architecture, by nesting an autoencoder inside another
 - The inner autoencoder is for temporal information
 - The outer autoencoder is for spatial information
- Spatial autoencoder is a simple undercomplete autoencoder
- Temporal autoencoder consists of a ConvLSTM as temporal encoder and optical flow as temporal decoder

Spatio-temporal Video Autoencoder

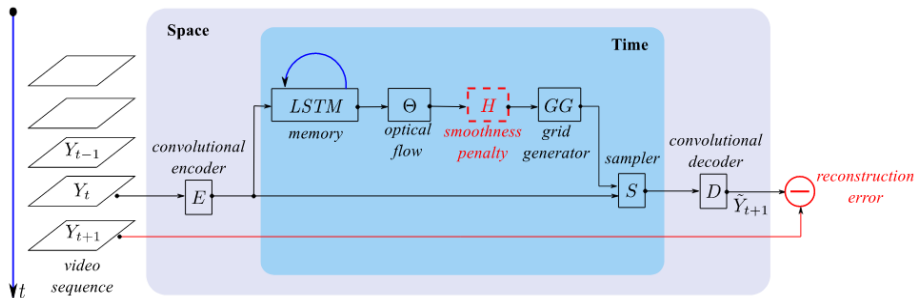


Figure: Spatio-temporal Video Autoencoder Architecture [11]

Spatio-temporal Video Autoencoder

- First input the whole image sequence X to the convolutional encoder (Works like a truncated SVD)

Spatio-temporal Video Autoencoder

- First input the whole image sequence X to the convolutional encoder (Works like a truncated SVD)
- Save the last output of the convolutional encoder

Spatio-temporal Video Autoencoder

- First input the whole image sequence X to the convolutional encoder (Works like a truncated SVD)
- Save the last output of the convolutional encoder
- Use the output sequence of the convolutional encoder as input for the ConvLSTM

Spatio-temporal Video Autoencoder

- First input the whole image sequence X to the convolutional encoder (Works like a truncated SVD)
- Save the last output of the convolutional encoder
- Use the output sequence of the convolutional encoder as input for the ConvLSTM
- Output of ConvLSTM will get into the optical flow convnet

Spatio-temporal Video Autoencoder

- First input the whole image sequence X to the convolutional encoder (Works like a truncated SVD)
- Save the last output of the convolutional encoder
- Use the output sequence of the convolutional encoder as input for the ConvLSTM
- Output of ConvLSTM will get into the optical flow convnet
- Normalize optical flow map and apply with saved image in image sampler

- First input the whole image sequence X to the convolutional encoder (Works like a truncated SVD)
- Save the last output of the convolutional encoder
- Use the output sequence of the convolutional encoder as input for the ConvLSTM
- Output of ConvLSTM will get into the optical flow convnet
- Normalize optical flow map and apply with saved image in image sampler
- „Resize“ next image.

- End-to-End differentiable

Spatio-temporal Video Autoencoder

- End-to-End differentiable
- Trained using BPTT with CE-Loss for MovingMNIST (Two digits per image)

Spatio-temporal Video Autoencoder

- End-to-End differentiable
- Trained using BPTT with CE-Loss for MovingMNIST (Two digits per image)
- Used for one-frame prediction (19 images as input, one as output)

- End-to-End differentiable
- Trained using BPTT with CE-Loss for MovingMNIST (Two digits per image)
- Used for one-frame prediction (19 images as input, one as output)
- Architecture could be used for multi-frame prediction with a feedback loop

Spatio-temporal Video Autoencoder



Figure: Results of MovingMNIST experiment [11]

- „Deep Predictive Coding Networks For Video Prediction And Unsupervised Learning“ by Lotter et al. [8]

- „Deep Predictive Coding Networks For Video Prediction And Unsupervised Learning“ by Lotter et al. [8]
- Informally named PredNet

- „Deep Predictive Coding Networks For Video Prediction And Unsupervised Learning“ by Lotter et al. [8]
- Informally named PredNet
- Architecture is based on the concept of „predictive coding“.

- „Deep Predictive Coding Networks For Video Prediction And Unsupervised Learning“ by Lotter et al. [8]
- Informally named PredNet
- Architecture is based on the concept of „predictive coding“.
 - Neuro-scientific task of „continually making predictions of incoming sensory stimuli“.

- „Deep Predictive Coding Networks For Video Prediction And Unsupervised Learning“ by Lotter et al. [8]
- Informally named PredNet
- Architecture is based on the concept of „predictive coding“.
 - Neuro-scientific task of „continually making predictions of incoming sensory stimuli“.
- Network consists of arbitrary amount of layers

- „Deep Predictive Coding Networks For Video Prediction And Unsupervised Learning“ by Lotter et al. [8]
- Informally named PredNet
- Architecture is based on the concept of „predictive coding“.
 - Neuro-scientific task of „continually making predictions of incoming sensory stimuli“.
- Network consists of arbitrary amount of layers
- Depth of network is hyperparameter

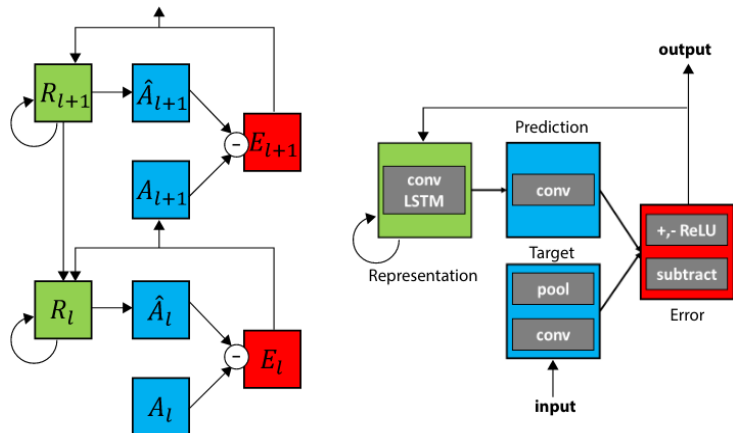


Figure: PredNet Architecture [8]

- $$A_l^t = \begin{cases} x_t & l = 0 \\ \text{MaxPool}(\text{ReLU}(\text{Conv}(E_{l-1}^t))) & l > 0 \end{cases} \quad (30)$$

- $$A_l^t = \begin{cases} x_t & l = 0 \\ \text{MaxPool}(\text{ReLU}(\text{Conv}(E_{l-1}^t))) & l > 0 \end{cases} \quad (30)$$

- $$\hat{A}_l^t = \text{ReLU}(\text{Conv}(R_l^t)) \quad (31)$$

- $$A_l^t = \begin{cases} x_t & l = 0 \\ \text{MaxPool}(\text{ReLU}(\text{Conv}(E_{l-1}^t))) & l > 0 \end{cases} \quad (30)$$

- $$\hat{A}_l^t = \text{ReLU}(\text{Conv}(R_l^t)) \quad (31)$$

- $$E_l^t = [\text{ReLU}(\hat{A}_l^t - A_l^t); \text{ReLU}(A_l^t - \hat{A}_l^t)] \quad (32)$$

- $$A_l^t = \begin{cases} x_t & l = 0 \\ \text{MaxPool}(\text{ReLU}(\text{Conv}(E_{l-1}^t))) & l > 0 \end{cases} \quad (30)$$

- $$\hat{A}_l^t = \text{ReLU}(\text{Conv}(R_l^t)) \quad (31)$$

- $$E_l^t = [\text{ReLU}(\hat{A}_l^t - A_l^t); \text{ReLU}(A_l^t - \hat{A}_l^t)] \quad (32)$$

- $$R_l^t = \text{ConvLSTM}(E_l^{t-1}, R_l^{t-1}, \text{Upsample}(R_{l+1}^t)) \quad (33)$$

- Doing prediction in every time-step instead after whole sequence

- Doing prediction in every time-step instead after whole sequence
- First time-step will output a spatially uniform prediction

- Doing prediction in every time-step instead after whole sequence
- First time-step will output a spatially uniform prediction
 - Because R_l and E_l are initialized with zero

- Doing prediction in every time-step instead after whole sequence
- First time-step will output a spatially uniform prediction
 - Because R_I and E_I are initialized with zero
- First compute R_I^t states in the top-down pass

- Doing prediction in every time-step instead after whole sequence
- First time-step will output a spatially uniform prediction
 - Because R_l and E_l are initialized with zero
- First compute R_l^t states in the top-down pass
- Then compute \hat{A}_l^t in forward pass

- Doing prediction in every time-step instead after whole sequence
- First time-step will output a spatially uniform prediction
 - Because R_l and E_l are initialized with zero
- First compute R_l^t states in the top-down pass
- Then compute \hat{A}_l^t in forward pass
- If current layer is *neg* lowest layer, compute A_l^t

- Doing prediction in every time-step instead after whole sequence
- First time-step will output a spatially uniform prediction
 - Because R_l and E_l are initialized with zero
- First compute R_l^t states in the top-down pass
- Then compute \hat{A}_l^t in forward pass
- If current layer is *neg* lowest layer, compute A_l^t
 - A_0^t is an image of the input sequence

- End-to-End differentiable

- End-to-End differentiable
- Trained using BPTT

- End-to-End differentiable
- Trained using BPTT
- Using an in-build error function

- End-to-End differentiable
- Trained using BPTT
- Using an in-build error function
 - Network consists of two modes (Error and Prediction mode)

- End-to-End differentiable
- Trained using BPTT
- Using an in-build error function
 - Network consists of two modes (Error and Prediction mode)

$$L_t = \sum_t \lambda_t \sum_l \frac{\lambda_l}{n_l} \sum_{n_l} E_l^t \quad (34)$$

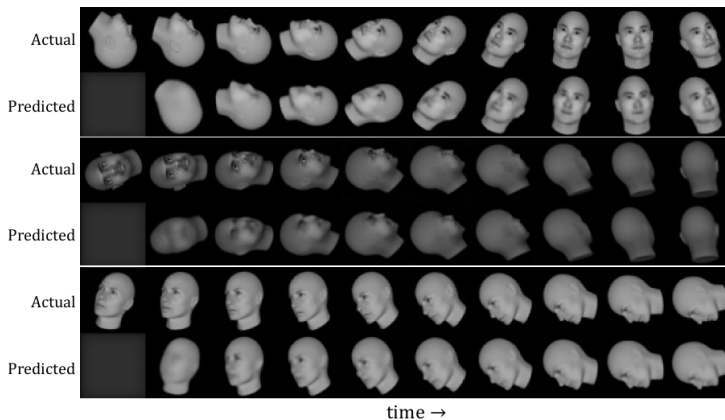
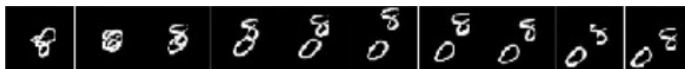
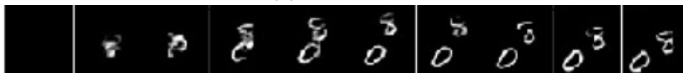


Figure: Results of FaceGen [8]



(a) Ground truth



(b) ConvLSTM

Figure: Results for MovingMNIST [3].

- „PredRNN: Recurrent Neural Networks for Predictive Learning using Spatiotemporal LSTMs“ by Wang et al. [16]

- „PredRNN: Recurrent Neural Networks for Predictive Learning using Spatiotemporal LSTMs“ by Wang et al. [16]
- Novel LSTM module named Spatiotemporal LSTM (ST-LSTM)

- „PredRNN: Recurrent Neural Networks for Predictive Learning using Spatiotemporal LSTMs“ by Wang et al. [16]
- Novel LSTM module named Spatiotemporal LSTM (ST-LSTM)
 - Consists of the typical temporal memory and extends this with an additional spatiotemporal memory

- „PredRNN: Recurrent Neural Networks for Predictive Learning using Spatiotemporal LSTMs“ by Wang et al. [16]
- Novel LSTM module named Spatiotemporal LSTM (ST-LSTM)
 - Consists of the typical temporal memory and extends this with an additional spatiotemporal memory
- Novel Architecture named PredRNN

$$g_t = \tanh(W_{xg} * X_t + W_{hg} * H_{t-1}^l + b_g) \quad (35)$$

$$i_t = \sigma(W_{xi} * X_t + W_{hi} * H_{t-1}^l + b_i) \quad (36)$$

$$f_t = \sigma(W_{xf} * X_t + W_{hf} * H_{t-1}^l + b_f) \quad (37)$$

$$C_t^l = f_t \odot C_{t-1}^l + i_t \odot g_t \quad (38)$$

$$g_t^l = \tanh(W_{xg'} * X_t + W_{mg} * M_t^{l-1} + b_{g'}) \quad (39)$$

$$i_t^l = \sigma(W_{xi'} * X_t + W_{mi} * M_t^{l-1} + b_{i'}) \quad (40)$$

$$f_t^l = \sigma(W_{xf'} * X_t + W_{mf} * M_t^{l-1} + b_{f'}) \quad (41)$$

$$M_t^l = f_t^l \odot M_t^{l-1} + i_t^l \odot g_t^l \quad (42)$$

$$o_t = \sigma(W_{xo} * X_t + W_{ho} * H_{t-1}^l + W_{co} * C_t^l + W_{mo} * M_t^l + b_o) \quad (43)$$

$$H_t^l = o_t \odot \tanh(W_{1 \times 1} * [C_t^l, M_t^l]) \quad (44)$$

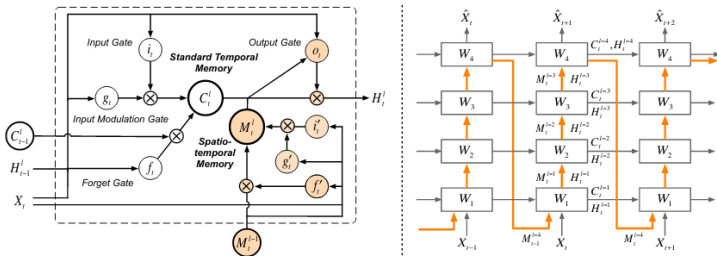


Figure: **Left:** ST-LSTM architecture, **Right:** PredRNN architecture [16]

- End-to-End differentiable

- End-to-End differentiable
- Trained using BPTT using CE-Loss on MovingMNIST (Two digits per frame)

- End-to-End differentiable
- Trained using BPTT using CE-Loss on MovingMNIST (Two digits per frame)
- Using ten images as input and output ten images

PredRNN



Figure: Results for MovingMNIST [16].

- Architectures have different approach but similar result

Related work

- Architectures have different approach but similar result
- First architecture uses LSTM in autoencoder architecture

- Architectures have different approach but similar result
- First architecture uses LSTM in autoencoder architecture
- The second architecture replaces the LSTM with ConvLSTM

- Architectures have different approach but similar result
- First architecture uses LSTM in autoencoder architecture
- The second architecture replaces the LSTM with ConvLSTM
- Third architecture leverages optical flow with two autoencoder, two split information in space and time

- Architectures have different approach but similar result
- First architecture uses LSTM in autoencoder architecture
- The second architecture replaces the LSTM with ConvLSTM
- Third architecture leverages optical flow with two autoencoder, two split information in space and time
- PredNet introduces a neuro-scientific idea to achieve state-of-the-art image-/video-prediction

- Architectures have different approach but similar result
- First architecture uses LSTM in autoencoder architecture
- The second architecture replaces the LSTM with ConvLSTM
- Third architecture leverages optical flow with two autoencoder, two split information in space and time
- PredNet introduces a neuro-scientific idea to achieve state-of-the-art image-/video-prediction
- The last architecture introduces a novel ConvLSTM named ST-LSTM and also a novel architecture named PredRNN

- Architectures have different approach but similar result
- First architecture uses LSTM in autoencoder architecture
- The second architecture replaces the LSTM with ConvLSTM
- Third architecture leverages optical flow with two autoencoder, two split information in space and time
- PredNet introduces a neuro-scientific idea to achieve state-of-the-art image-/video-prediction
- The last architecture introduces a novel ConvLSTM named ST-LSTM and also a novel architecture named PredRNN
- \Rightarrow Even though, the algorithms have different approaches to perform state-of-the-art image-/video-prediction, they all use some kind of LSTM as recurrent module.

Contents

- 1 Scientific questions
- 2 Introduction
- 3 Machine Learning Theory
 - Image Prediction
 - Autoencoder
 - CNN
 - RNN
 - LSTM
 - ConvLSTM
 - Backpropagation
 - BPTT
- 4 Image Prediction Architectures
 - LSTM Autoencoder
 - ConvLSTM Autoencoder
 - Spatio-temporal Video Autoencoder
 - PredNet
 - PredRNN
- 5 Experiments
 - Experimental setup
 - Experimental results (First setup)

Experimental setup

- Use the three implemented networks and change the recurrent layer with PredRNN to see if this achieves better performance

Experimental setup

- Use the three implemented networks and change the recurrent layer with PredRNN to see if this achieves better performance
- Experiments were performed on three different datasets

Experimental setup

- Use the three implemented networks and change the recurrent layer with PredRNN to see if this achieves better performance
- Experiments were performed on three different datasets
 - MovingMNIST (synthetic dataset)

Experimental setup

- Use the three implemented networks and change the recurrent layer with PredRNN to see if this achieves better performance
- Experiments were performed on three different datasets
 - MovingMNIST (synthetic dataset)
 - Pre-processed to have frame size $(1 \times 64 \times 64)$, two digits per frame

Experimental setup

- Use the three implemented networks and change the recurrent layer with PredRNN to see if this achieves better performance
- Experiments were performed on three different datasets
 - MovingMNIST (synthetic dataset)
 - Pre-processed to have frame size ($1 \times 64 \times 64$), two digits per frame
 - Training set of 10000 sequences, with every sequence consisting of ten images

Experimental setup

- Use the three implemented networks and change the recurrent layer with PredRNN to see if this achieves better performance
- Experiments were performed on three different datasets
 - MovingMNIST (synthetic dataset)
 - Pre-processed to have frame size ($1 \times 64 \times 64$), two digits per frame
 - Training set of 10000 sequences, with every sequence consisting of ten images
 - Test set of 3000 sequences

Experimental setup

- Use the three implemented networks and change the recurrent layer with PredRNN to see if this achieves better performance
- Experiments were performed on three different datasets
 - MovingMNIST (synthetic dataset)
 - Pre-processed to have frame size ($1 \times 64 \times 64$), two digits per frame
 - Training set of 10000 sequences, with every sequence consisting of ten images
 - Test set of 3000 sequences
 - KTH (real dataset, fixed camera)

Experimental setup

- Use the three implemented networks and change the recurrent layer with PredRNN to see if this achieves better performance
- Experiments were performed on three different datasets
 - MovingMNIST (synthetic dataset)
 - Pre-processed to have frame size ($1 \times 64 \times 64$), two digits per frame
 - Training set of 10000 sequences, with every sequence consisting of ten images
 - Test set of 3000 sequences
 - KTH (real dataset, fixed camera)
 - Pre-processed to gray-scaled images, cropped to size ($1 \times 80 \times 60$)

Experimental setup

- Use the three implemented networks and change the recurrent layer with PredRNN to see if this achieves better performance
- Experiments were performed on three different datasets
 - MovingMNIST (synthetic dataset)
 - Pre-processed to have frame size ($1 \times 64 \times 64$), two digits per frame
 - Training set of 10000 sequences, with every sequence consisting of ten images
 - Test set of 3000 sequences
 - KTH (real dataset, fixed camera)
 - Pre-processed to gray-scaled images, cropped to size ($1 \times 80 \times 60$)
 - Kitti (real dataset, ego-centric movement)

Experimental setup

- Use the three implemented networks and change the recurrent layer with PredRNN to see if this achieves better performance
- Experiments were performed on three different datasets
 - MovingMNIST (synthetic dataset)
 - Pre-processed to have frame size ($1 \times 64 \times 64$), two digits per frame
 - Training set of 10000 sequences, with every sequence consisting of ten images
 - Test set of 3000 sequences
 - KTH (real dataset, fixed camera)
 - Pre-processed to gray-scaled images, cropped to size ($1 \times 80 \times 60$)
 - Kitti (real dataset, ego-centric movement)
 - Cropped to frame size ($3 \times 160 \times 128$)

Experimental setup

- First experiments were performed using fixed hyperparameter (HP optimization is only applied in the second experiments)

Experimental setup

- First experiments were performed using fixed hyperparameter (HP optimization is only applied in the second experiments)
 - In total 18 different networks trained

Experimental setup

- First experiments were performed using fixed hyperparameter (HP optimization is only applied in the second experiments)
 - In total 18 different networks trained
 - On own Computer \Rightarrow Kitti training took $> 12h$

Experimental setup

- First experiments were performed using fixed hyperparameter (HP optimization is only applied in the second experiments)
 - In total 18 different networks trained
 - On own Computer \Rightarrow Kitti training took $> 12h$
- MovingMNIST 5000 iterations per epoch in one epoch (Inspired by Elsayed et al. [3])

Experimental setup

- First experiments were performed using fixed hyperparameter (HP optimization is only applied in the second experiments)
 - In total 18 different networks trained
 - On own Computer \Rightarrow Kitti training took $> 12h$
- MovingMNIST 5000 iterations per epoch in one epoch (Inspired by Elsayed et al. [3])
- KTH 500 iterations per epoch in 20 epochs

Experimental setup

- First experiments were performed using fixed hyperparameter (HP optimization is only applied in the second experiments)
 - In total 18 different networks trained
 - On own Computer \Rightarrow Kitti training took $> 12h$
- MovingMNIST 5000 iterations per epoch in one epoch (Inspired by Elsayed et al. [3])
- KTH 500 iterations per epoch in 20 epochs
- Kitti 500 iterations per epoch in 50 epochs

Experimental setup

- First experiments were performed using fixed hyperparameter (HP optimization is only applied in the second experiments)
 - In total 18 different networks trained
 - On own Computer \Rightarrow Kitti training took $> 12h$
- MovingMNIST 5000 iterations per epoch in one epoch (Inspired by Elsayed et al. [3])
- KTH 500 iterations per epoch in 20 epochs
- Kitti 500 iterations per epoch in 50 epochs
- Adam [6] was used as optimizer with lr of 0.001 and lr-scheduler reducing the lr after $\frac{epochs}{2}$ by factor ten

Experimental setup

- First experiments were performed using fixed hyperparameter (HP optimization is only applied in the second experiments)
 - In total 18 different networks trained
 - On own Computer \Rightarrow Kitti training took $> 12\text{h}$
- MovingMNIST 5000 iterations per epoch in one epoch (Inspired by Elsayed et al. [3])
- KTH 500 iterations per epoch in 20 epochs
- Kitti 500 iterations per epoch in 50 epochs
- Adam [6] was used as optimizer with lr of 0.001 and lr-scheduler reducing the lr after $\frac{\text{epochs}}{2}$ by factor ten
- Other hyperparameter can be found in the yml-folder of the implementation

- Second experiments are performed using HP optimization, using early-stopping and with „optimal values“

- Second experiments are performed using HP optimization, using early-stopping and with „optimal values“
- Only one out of 18 experiments performed, to get the importance of HP optimization and a tendency of how the results should look using optimal conditions

Experimental results (First setup)

- Section is splitted into the three different datasets

Experimental results (First setup)

- Section is splitted into the three different datasets
- First comparison of trainable parameter amount

Experimental results (First setup)

- Section is splitted into the three different datasets
- First comparison of trainable parameter amount
- Secondly table of mean MSE for testing and some example output

Model	ConvLSTM	PredRNN
Autoencoder (Depth 2)	537.411	773.018
PredNet	6.909.818	12.421.090
Spatiotemp	1.001.324	1.415.639

Table: Number of trainable parameter for MovingMNIST.

MovingMNIST

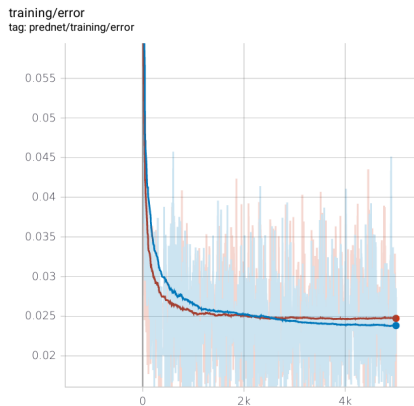
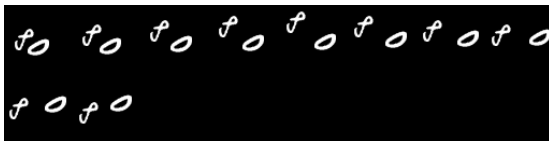
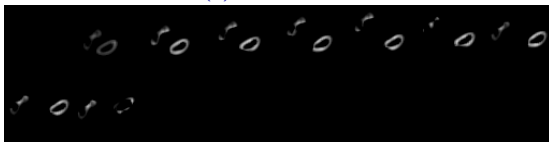


Figure: Training PredNet with ConvLSTM (blue) and PredRNN (red) on MovingMNIST.

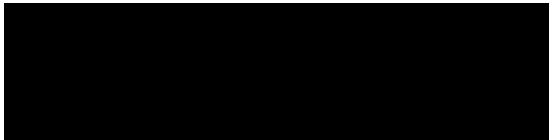
MovingMNIST



(a) Ground truth



(b) ConvLSTM



(c) PredRNN (vanished)

Figure: Test results for PredNet on MovingMNIST.

Model	ConvLSTM	PredRNN
Autoencoder (Depth 2)	0.027	0.028
PredNet	0.035	0.041
Spatiotemp	0.024	0.022

Table: Mean MSE for MovingMNIST.

Model	ConvLSTM	PredRNN
Autoencoder (Depth 2)	542.531	781.760
PredNet	850.325	1.285.730
Spatiotemp	1.007.598	1.421.913

Table: Number of trainable parameter for KTH.



(a) Ground truth



(b) ConvLSTM



(c) PredRNN

Figure: Test results for PredNet on KTH.

Model	ConvLSTM	PredRNN
Autoencoder (Depth 2)	$1.55e - 3$	0.05 (didn't converged)
PredNet	$1.95e - 3$	$1.93e - 3$
Spatiotemp	$3.1e - 3$	0.025 (didn't converged)

Table: Mean MSE for KTH.

Model	ConvLSTM	PredRNN
Autoencoder (Depth 2)	542.531	781.760
PredNet	8.222.559	12.430.626
Spatiotemp	1.640.321	2.200.641

Table: Number of trainable parameter for Kitti.



(a) Ground truth



(b) ConvLSTM



(c) PredRNN

Figure: Test results for PredNet on Kitti.

Model	ConvLSTM	PredRNN
Autoencoder (Depth 2)	0.02	0.013
PredNet	0.019	0.02
Spatiotemp	0.018	0.017

Table: Mean MSE for Kitti.

Experimental Results (Second setup)

- Performed on KTH dataset with optimal values and lr of 0.0001

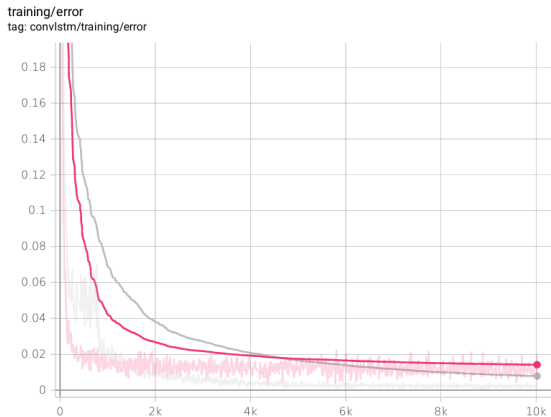


Figure: Training Autoencoder with ConvLSTM (red) and PredRNN (gray) on KTH using hyperparameter optimization and early stopping.

Experimental Results (Second setup)

- Mean test MSE for Autoencoder using ConvLSTM: 0.0013

Experimental Results (Second setup)

- Mean test MSE for Autoencoder using ConvLSTM: 0.0013
- Mean test MSE for Autoencoder using PredRNN: 0.00072

Experimental Results (Second setup)

- Mean test MSE for Autoencoder using ConvLSTM: 0.0013
- Mean test MSE for Autoencoder using PredRNN: 0.00072
- PredRNN is able to boost Autoencoder by more then 80 percent

Experimental Results (Second setup)

- Mean test MSE for Autoencoder using ConvLSTM: 0.0013
- Mean test MSE for Autoencoder using PredRNN: 0.00072
- PredRNN is able to boost Autoencoder by more then 80 percent
- Showed importance of HP optimization and early stopping

Experimental Results (Second setup)

- Mean test MSE for Autoencoder using ConvLSTM: 0.0013
- Mean test MSE for Autoencoder using PredRNN: 0.00072
- PredRNN is able to boost Autoencoder by more then 80 percent
- Showed importance of HP optimization and early stopping
- Gives a hint of how the other experiments should perform

Experimental Results (Second setup)



(a) Ground truth



(b) Con-
vLSTM



(c)
PredRNN

Figure: Test results for Autoencoder on KTH using hyperparameter optimization and early stopping.

Contents

- 1 Scientific questions
- 2 Introduction
- 3 Machine Learning Theory
 - Image Prediction
 - Autoencoder
 - CNN
 - RNN
 - LSTM
 - ConvLSTM
 - Backpropagation
 - BPTT
- 4 Image Prediction Architectures
 - LSTM Autoencoder
 - ConvLSTM Autoencoder
 - Spatio-temporal Video Autoencoder
 - PredNet
 - PredRNN
- 5 Experiments
 - Experimental setup
 - Experimental results (First setup)

- Related work showed different state-of-the-art solutions for image-/video-prediction

- Related work showed different state-of-the-art solutions for image-/video-prediction
 - Showed how previous work is used and adapted to gain better performance

- Related work showed different state-of-the-art solutions for image-/video-prediction
 - Showed how previous work is used and adapted to gain better performance
 - E.g. Shi et al. [14] used the autoencoder architecture from Srivastava et al. [15] but with novel ConvLSTM module

- Related work showed different state-of-the-art solutions for image-/video-prediction
 - Showed how previous work is used and adapted to gain better performance
 - E.g. Shi et al. [14] used the autoencoder architecture from Srivastava et al. [15] but with novel ConvLSTM module
 - All have a different approach to get to the same goal and rely on the same recurrent module basis

- Related work showed different state-of-the-art solutions for image-/video-prediction
 - Showed how previous work is used and adapted to gain better performance
 - E.g. Shi et al. [14] used the autoencoder architecture from Srivastava et al. [15] but with novel ConvLSTM module
 - All have a different approach to get to the same goal and rely on the same recurrent module basis
- Network implementation showed, that a researcher should have basic skills in Python and more then one deep learning framework, in order to be able to read and understand code

- Related work showed different state-of-the-art solutions for image-/video-prediction
 - Showed how previous work is used and adapted to gain better performance
 - E.g. Shi et al. [14] used the autoencoder architecture from Srivastava et al. [15] but with novel ConvLSTM module
 - All have a different approach to get to the same goal and rely on the same recurrent module basis
- Network implementation showed, that a researcher should have basic skills in Python and more then one deep learning framework, in order to be able to read and understand code
- First experiments were fixed and didn't showed any superior behavior of PredRNN as recurrent layer

- Related work showed different state-of-the-art solutions for image-/video-prediction
 - Showed how previous work is used and adapted to gain better performance
 - E.g. Shi et al. [14] used the autoencoder architecture from Srivastava et al. [15] but with novel ConvLSTM module
 - All have a different approach to get to the same goal and rely on the same recurrent module basis
- Network implementation showed, that a researcher should have basic skills in Python and more then one deep learning framework, in order to be able to read and understand code
- First experiments were fixed and didn't showed any superior behavior of PredRNN as recurrent layer
- Using HP optimization, the tables turned and even one example outperformed the ConvLSTM solution with a performance boost of $> 80\%$

Contents

- 1 Scientific questions
- 2 Introduction
- 3 Machine Learning Theory
 - Image Prediction
 - Autoencoder
 - CNN
 - RNN
 - LSTM
 - ConvLSTM
 - Backpropagation
 - BPTT
- 4 Image Prediction Architectures
 - LSTM Autoencoder
 - ConvLSTM Autoencoder
 - Spatio-temporal Video Autoencoder
 - PredNet
 - PredRNN
- 5 Experiments
 - Experimental setup
 - Experimental results (First setup)

Conclusion

- Different approaches for image-/video-prediction were discussed

Conclusion

- Different approaches for image-/video-prediction were discussed
- The choice of the recurrent module is important for the performance of the network

- Different approaches for image-/video-prediction were discussed
- The choice of the recurrent module is important for the performance of the network
- The results showed, how important HP optimization is to achieve state-of-the-art performance and that without, having a more complex sub-module is not the key

Conclusion

- Different approaches for image-/video-prediction were discussed
- The choice of the recurrent module is important for the performance of the network
- The results showed, how important HP optimization is to achieve state-of-the-art performance and that without, having a more complex sub-module is not the key
- For future experiments, it would be nice to have all 18 experiments performed with HP optimization to get the overall results

Bibliography

- [1] CNN — Introduction to Pooling Layer. <https://www.geeksforgeeks.org/cnn-introduction-to-pooling-layer/>. – Accessed: 2020-08-27
- [2] Understanding LSTM Networks. <https://colah.github.io/posts/2015-08-Understanding-LSTMs/>. – Accessed: 2020-07-13
- [3] ELSAYED, Nelly ; MAIDA, Anthony S. ; BAYOUMI, Magdy A.: Reduced-Gate Convolutional LSTM Using Predictive Coding for Spatiotemporal Prediction. In: *CoRR* abs/1810.07251 (2018). – URL <http://arxiv.org/abs/1810.07251>
- [4] GOODFELLOW, Ian ; BENGIO, Yoshua ; COURVILLE, Aaron: *Deep Learning*. MIT Press, 2016. – <http://www.deeplearningbook.org>
- [5] HOCHREITER, Sepp ; SCHMIDHUBER, Jürgen: Long Short-term Memory. In: *Neural computation* 9 (1997), 12, S. 1735–80
- [6] KINGMA, Diederik P. ; BA, Jimmy: Adam: A Method for Stochastic Optimization. In: *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, URL <http://arxiv.org/abs/1412.6980>, 2015
- [7] LECUN, Y. ; BOTTOU, L. ; BENGIO, Y. ; HAFNER, P.: Gradient-based