

TSSL Lab 4 - Recurrent Neural Networks

In this lab we will explore different RNN models and training procedures for a problem in time series prediction.

```
In [1]: import numpy as np
import tensorflow as tf
from tensorflow.keras import keras
from tensorflow.keras import layers
import tensorflow.keras.backend as K
import matplotlib.pyplot as plt
from sklearn.metrics import mean_squared_error as mse

plt.rcParams['figure.figsize'] = (10,6) # Increase default size of plots

Set the random seed for reproducibility

In [2]: np.random.seed(42)
tf.random.set_seed(42)
```

1. Load and prepare the data

We will build a model for predicting the number of **sunspots**. We work with a data set that has been published on [Kaggle](#), with the description:

Sunspots are rotational phenomena on the Sun's photosphere that appear as spots darker than the surrounding areas. They are regions of reduced surface temperature caused by concentrations of magnetic field flux that inhibit convection. Sunspots usually appear in pairs of opposite magnetic polarity. Their number varies according to the approximately 11-year solar cycle.

The data consists of the monthly mean total sunspot number, from 1749-01-01 to 2017-08-31.

```
In [3]: # Read the data
data = pandas.read_csv('sunspots.csv', header=0)
dates = data['Date'].values
y = data['Monthly Mean Total Sunspot Number'].values
ndata = len(y)
print(f'Total number of data points: {ndata}')

# We define a train/test split, here with 70 % training data
ntrain = int(ndata*0.7)
ntest = ndata-ntrain
print(f'Number of training data points: {ntrain}')
Total number of data points: 3252
Number of training data points: 2276

In [4]: plt.plot(dates[ntrain:], y[ntrain:])
plt.plot(dates[ntrain:], y[ntrain:])
plt.xticks(range(0, ndata, 300), dates[1:300], rotation = 90); # Show only one tick every 25th year for clarity
```

There is a clear seasonality to the data, but the amplitude of the peaks varies quite a lot. Also, we note that the data is nonnegative, which is natural since it consists of counts of sunspots. However, for simplicity we will not take this constraint into account in this lab assignment and allow ourselves to model the data using a Gaussian likelihood (i.e. using MSE as a loss function).

From the plot we see that the range of the data is roughly [0,400] so as a simple normalization we divide by the constant `MAX_VAL=400`.

```
In [5]: MAX_VAL = 400
y = y/MAX_VAL
```

2. Baseline methods

Before constructing any sophisticated models using RNNs, let's consider two baseline methods,

1. The first baseline is a "naive" method which simply predicts $y_t = y_{t-1}$.
2. The second baseline is an AR(p) model (based on the implementation used for lab 1).

We evaluate the performance of these methods in terms of mean-squared-error and mean-absolute-error, to compare the more advanced models with later on.

```
In [6]: def evaluate_performance(y_pred, y, split_time, name=None):
    """This function evaluates and prints the MSE and MAE of the prediction.

    Parameters
    -----
    y_pred: ndarray
        Array of size (n,) with predictions.
    y: ndarray
        Array of size (n,) with target values.
    split_time: int
        The leading number of elements in y_pred and y that belong to the training data set.
        The remaining elements, i.e. y_pred[split_time:] and y[split_time:] are treated as test data.
    """

    # Compute error in prediction
    train_mse = mse(y_pred[:split_time], y[:split_time])
    test_mse = mse(y_pred[split_time:], y[split_time:])

    # We evaluate the MSE and MAE in the original scale of the data, i.e. we add back MAX_VAL
    train_mse = np.mean(resid(split_time**2)*MAX_VAL**2)
    test_mse = np.mean(resid(split_time**2)*MAX_VAL**2)
    train_mae = np.mean(np.abs(resid(split_time))*MAX_VAL)
    test_mae = np.mean(np.abs(resid(split_time))*MAX_VAL)

    # Print
    print(f'Model {name}\n Training MSE: {train_mse:.4f}, MAE: {train_mae:.4f}\n Testing MSE: {test_mse:.4f}, MAE: {test_mae:.4f}')

In [7]: # Store the predictions in an array of length ndata-1. Note that there is a shift in the indices
# between the prediction and the observation sequence, since there is no prediction available for the first observation.
# Specifically, y_pred_naive(t) is a prediction of y(t+1), so the first element of y_pred_naive is a prediction of the second element of y, and so on. We will use the same "bookkeeping convention" throughout the lab, so it is in
# you understand it.

y_pred_naive = y[ntrain:]

evaluate_performance(y_pred_naive, # Predictions
                    y[1:], # Corresponding target values
                    ntrain, # Number of leading elements in the input arrays corresponding to training data
                    name='Naive')

Model Naive
Training MSE: 776.3437, MAE: 19.3285
Testing MSE: 708.0360, MAE: 19.2256
```

Next, we consider a slightly more advanced baseline method, namely an AR(p) model.

```
In [8]: # We import two functions that were written as part of lab 1
from tsstools.lab4 import fit_ar, predict_ar_step

# Order of the AR model (set by a few manual trials)
ar_coef = fit_ar(y[ntrain:], p) # Fit the model to the training data

# Predict. Note that y contains both training and validation data,
# and the prediction is for the values y[1:], ..., y(n).
y_pred_ar = predict_ar_step(ar_coef, y)

In [9]: evaluate_performance(y_pred_ar, # The prediction array is of length n-p
                        y[1:], # Corresponding target values
                        ntrain, # Number of leading elements in the input arrays corresponding to training data
                        name='AR')

Model AR
Training MSE: 603.8656, MAE: 17.3420
Testing MSE: 590.3752, MAE: 17.6221
```

3. Simple RNN

We will now construct a model based on a recurrent neural network. We will initially use the `SimpleRNN` class from `Keras`, which correspond to the basic Jordan-Elman network presented in the lectures.

Q2: Assume that we construct an "RNN cell" using the call `layers.SimpleRNN(units = d, return_sequences=True)`. Now, assume that an array `X` with the dimensions `[Q,M,P]` is fed as the input to the above object. We know that `X` contains a set of sequences (time series) with equal lengths. Specify which of the symbols `Q, M, P` that corresponds to each of the items below.

- The length of the sequences (number of time steps)
- The number of features (at each time step), i.e. the dimension of each time series
- The number of sequences

Furthermore, specify the values of Q, M, P for the data at hand (treated as a single time series).

Hint: Read the documentation for `SimpleRNN` to find the answer.

A2: `Q` is the number of sequences, `M` is the number of features, and `P` is the number of time steps. `Q` is the number of sequences, `M` is the number of features, and `P` is the number of time steps. `Q` is the number of sequences, `M` is the number of features, and `P` is the number of time steps. `Q` is the number of sequences, `M` is the number of features, and `P` is the number of time steps.

Q3: Continuing the question above, answer the following:

- What is the meaning of setting `units = d`?
- Assume that we pass a single time series of length n as input to the layer. Then what is the dimension of the output?
- If we would had set the parameter `return_sequences=False` when constructing the layer, then what would be the answer to the previous question?

A3: Units are the positive integer, dimensionality of the output space.

There will be two output: whose sequence of $(n,1,d)$ and final state will have (n,d)

The output of `return_sequence = False`, then the output will be (n,d) dimension

In `Keras`, each layer is created separately and are then joined by a `Sequential` object. It is very easy to construct stacked models in this way. The code below corresponds to a simple Jordan-Elman Network on the form,

$$\hat{y}_{t+1} = \sigma(W_h \cdot h_t + U y_t + b),$$
$$\hat{y}_{t+1} = C h_t + c,$$

Note: It is not necessary to explicitly specify the input shape, since this can be inferred from the input on the first call. However, for the `summary` function to work we need to tell the model what the dimension of the input is so that it can infer the correct sizes of the involved matrices. Also note that in `Keras` you can sometimes use `None` when some dimensions are not known in advance.

```
In [10]: d = 10 # hidden state dimension

model1 = keras.Sequential([
    # Simple RNN layer
    layers.SimpleRNN(units = d, input_shape=(None,1), return_sequences=True, activation='tanh'),
    # A linear output layer
    layers.Dense(units = 1, activation='linear')
])

# We store the initial weights in order to get an exact copy of the model when trying different training procedures
init_weights = model1.get_weights().copy()
Model: "sequential"

Layer (type) Output Shape Param #
=====
simple_rnn (SimpleRNN) (None, None, 10) 110
dense (Dense) (None, None, 1) 1
-----
Total params: 131
Trainable params: 131
Non-trainable params: 0
```

Q4: From the model summary we can see the number of parameters associated with each layer. Relate these numbers to the dimensions of the weight matrices and bias vectors $\{W, U, b, C, c\}$ in the mathematical model definition above.

A4: W will take parameters from the dense layer for each state dimension, it is $[11,10]$. U will have d rows, hence it is $[10,1]$. b is d rows, hence it is $[10,1]$. c is d rows, hence it is $[10,1]$. U will have d rows, hence it is $[10,1]$. b is d rows, hence it is $[10,1]$. c is d rows, hence it is $[10,1]$. U will have d rows, hence it is $[10,1]$. b is d rows, hence it is $[10,1]$. c is d rows, hence it is $[10,1]$.

4. Training the RNN model

In this section we will consider a few different ways of handling the data when training the simple RNN model constructed above. As a first step, however, we construct explicit input and target (output) arrays for the training and test data, which will simplify the calls to the training procedures below.

The task that we consider in this lab is one-step prediction, i.e. at each time step we compute a prediction $\hat{y}_{t+1} \approx y_t$ which depend on the previous observations $y_{1:t}$. However, when working with RNNs, the information contained in previous observations is aggregated in the state of the RNN, and we will only use y_{t-1} as the *explicit input* at each time step t .

Furthermore, when addressing a problem of time series prediction, it is often a good idea to introduce an explicit skip connection from the input y_{t-1} to the prediction y_t . Equivalently, we can define the *target value* at time step t to be the residual $\hat{y}_t := y_t - y_{t-1}$. Indeed, if the model can predict the value of the residual, then we can simply add back y_{t-1} to get a prediction of y_t .

Taking this into consideration, we define explicit input and output arrays as shifted versions of the data series $y_{1:n}$.

```
In [11]: # Training data
x_train = y[ntrain-1:] # Input is denoted by x, training inputs are x[0]=y[0], ..., x[ntrain-1]=y[ntrain-1]
y_train = y[ntrain:] # Output is denoted by y, training outputs are y[0]=y[1]-y[0], ..., y[ntrain-1]=y[ntrain]-y[ntrain-1]

# Test data
x_test = y[ntrain-1:] # Test inputs are x_test[0] = y[ntrain-1], ..., x_test[ntrain-1] = y[ntrain-1]
y_test = y[ntrain:] # Test outputs are y_test[0] = y[ntrain]-y[ntrain-1], ..., y_test[ntrain-1] = y[ntrain]-y[ntrain-1]

# Reshape the data
x_train = x_train.reshape((1,ntrain-1,1))
y_train = y_train.reshape((1,ntrain-1,1))
x_test = x_test.reshape((1,ntrain-1,1))
y_test = y_test.reshape((1,ntrain-1,1))

Option 1. Process all data in each gradient computation ("do nothing")

The first option is to process all data at each iteration of the gradient descent method.
```

```
In [12]: model1 = keras.models.clone_model(model0) # This creates a new instance of the same model
model1.set_weights(init_weights) # We set the initial weights to be the same for all models

# Evaluate MSE and MAE (both training and test data)
evaluate_performance(y_pred1, y[1:], ntrain-1, name='Simple RNN, "do nothing"')

Model Simple RNN, "do nothing"
Training MSE: 617.7963, MAE: 17.8611
Testing MSE: 573.7251, MAE: 17.7216

Predictions on test data
```

Q5: What should we set the `batch_size` to, in order to compute the gradient based on the complete training data sequence at each iteration? Complete the code below.

Note: You can set `verbose=1` if you want to monitor the training progress, but if you do, please **clear the output of the cell** before generating a pdf with your solutions, so that we don't get multiple pages with training errors in the submitted reports.

```
In [13]: model1.compile(loss='mse', optimizer='rmsprop', metrics=['mse'])
history = model1.fit(x_train, y_train,
                    epochs = 200,
                    batch_size = 1,
                    verbose = 0,
                    validation_data = (x_test, y_test))

We plot the training and test error vs the iteration (epoch) number, using a helper function from the tsstools_lab4 module.
```

```
In [14]: from tsstools.lab4 import plot_history
start_at = 10 # Skip the first few epochs for clarity
plot_history(history, start_at)
```

Q6: Finally we compute the predictions of $\{y_t\}$ for both the training and test data using the model's `predict` function. Complete the code below to compute the predictions.

Hint: You need to reshape the data when passing it to the `predict` to comply with the input shape used in `Keras` (cf. above).

Hint: Since the model is trained on the residuals \hat{y}_t , don't forget to add back y_{t-1} when predicting y_t . However, make sure that you don't "cheat" by using a non-causal predictor (i.e. using y_t when predicting y_t).

```
In [15]: # Predict on all data using the final model.
y_n = y.reshape((1,ndata-1,1))
# We predict using y_1,...,y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n.
y_pred1 = model1.predict(y_n) # We predict using y_1,...,y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n.
y_pred1 = model1.predict(y_n).flatten() + y[ntrain-1:]

Using the prediction computed above we can plot them and evaluate the performance of the model in terms of MSE and MAE.
```

```
In [16]: def plot_prediction(y_pred):
    # Plot prediction on test data
    plt.plot(dates[ntrain:], y_pred[ntrain-1:])
    # Plot prediction on training data
    plt.plot(dates[ntrain:], y_pred[ntrain-1:])
    plt.xticks(range(0, ntest, 300), dates[ntrain:300], rotation = 90); # Show only one tick every 25th year
    plt.legend(['Data', 'Prediction'])
    plt.title('Predictions on test data')

In [17]: # Plot prediction
plot_prediction(y_pred1)

# Evaluate MSE and MAE (both training and test data)
evaluate_performance(y_pred1, y[1:], ntrain-1, name='Simple RNN, "do nothing"')

Model Simple RNN, "do nothing"
Training MSE: 617.7963, MAE: 17.8611
Testing MSE: 573.7251, MAE: 17.7216

Predictions on test data
```

Option 2. Random windowing

Instead of using all the training data when computing the gradient for the numerical optimizer, we can speed it up by restricting the gradient computation to a smaller window of consecutive time steps. Here, we sample a random window within the training data and "pretend" that this window is independent from the observations outside the window. Specifically, when processing the observations within each window the hidden state of the RNN is initialized to zero at the first time point in the window.

To implement this method in Python, we will make use of a *generator function*. A generator is a function that can be paused, return an intermediate value, and then resumed to continue its execution. An intermediate return value is produced using the `yield` keyword. Generators are used in `Keras` to implement infinite loops that feed the training procedure with training data. Specifically, the `yield` statement of the generator should return a pair x_t, y_t with inputs and corresponding targets from the training data. Each epoch of the training procedure will then call the generator for a total of `steps_per_epoch` such `yield` statements.

```
In [18]: def generator_train(window_size):
    while True:
        # We sample a random window of size window_size from the training data.
        # Hence, the maximum index included in randint, so the maximum value that we can get is tt = ntrain-window_size.
        # start of window = np.random.randint(0, ntrain - window_size) # First time index of window (inclusive)
        end_of_window = start_of_window + window_size # Last time index of window (exclusive, i.e. this is not included)
        yield x_train[start_of_window:end_of_window, :], y_train[start_of_window:end_of_window, :]

In [19]: model2 = keras.models.clone_model(model0) # This creates a new instance of the same model
model2.set_weights(init_weights) # We set the initial weights to be the same for all models

# Evaluate MSE and MAE (both training and test data)
evaluate_performance(y_pred2, y[1:], ntrain-1, name='Simple RNN, windowing')

Model Simple RNN, windowing
Training MSE: 601.0721, MAE: 17.3281
Testing MSE: 572.9989, MAE: 17.1612

Predictions on test data
```

Q7: Assume that we process a window of observations of length `window_size` at each iteration. Then, how many gradient steps per epoch can we afford, for computational cost per epoch to be comparable to the method considered in Option 1? Set the `steps_per_epoch` parameter of the fitting function to your answer.

```
In [20]: window_size = 100
model2.compile(loss='mse', optimizer='rmsprop', metrics=['mse'])
history = model2.fit(generator_train(window_size),
                    epochs = 200,
                    batch_size = 1,
                    verbose = 0,
                    validation_data = (x_test, y_test),
                    steps_per_epoch = int(ndata/window_size))

Similarly to above we plot the error curves vs the iteration (epoch) number.
```

```
In [21]: plot_history(history, start_at)
```

Q8: Comparing this error plot to the one you got for training Option 1, can you see any *qualitative* differences? Explain the reason for the difference.

A8: The errors in Option 1 sunk continuously and converged around 0.004. Here in Option 2 the training error varies a lot (between 0.005 and 0.0025), yet the test error varies less and has a slight downward trend. This is because we fit the epochs to random windows, so depending on the window the error is higher or lower. Both errors in Option 2 seem slightly lower than in Option 1 because we do more steps per epoch and therefore reach a better solution faster.

Q9: Compute a prediction for all values of $\{y_2, \dots, y_n\}$ analogously to **Q6**.

```
In [22]: # Predict on all data using the final model.
y_n = y.reshape((1,ndata-1,1))
# We predict using y_1,...,y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n.
y_pred2 = model2.predict(y_n) # We predict using y_1,...,y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n.
y_pred2 = model2.predict(y_n).flatten() + y[ntrain-1:]

In [23]: # Plot prediction on test data
plot_prediction(y_pred2)

# Evaluate MSE and MAE (both training and test data)
evaluate_performance(y_pred2, y[1:], ntrain-1, name='Simple RNN, windowing')

Model Simple RNN, windowing
Training MSE: 601.0721, MAE: 17.3281
Testing MSE: 572.9989, MAE: 17.1612

Predictions on test data
```

Option 3. Sequential windowing with stateful training

As a final option, we consider a model aimed at better respecting the temporal dependencies between calls. This is based on "statefulness" which simply means that the RNN remembers its hidden state between calls. That is, if model is in stateful mode and is used to process two sequences of inputs after each other, then the final state from the first sequence is used as the initial state for the second sequence.

```
In [24]: # To enable stateful training, we need to create model where we set stateful=True in the RNN layer
model3 = keras.Sequential([
    # Simple RNN layer with stateful=True
    layers.SimpleRNN(units = d, batch_input_shape=(1, None, 1), return_sequences=True, stateful=True, activation='tanh'),
    # A linear output layer
    layers.Dense(units = 1, activation='linear')
])
model3.set_weights(init_weights)

Q10: When working with stateful training we need to make some adjustments to the training data generator.
```

1. First, the RNN model doesn't keep track of the actual time indices of the different windows that it is fed. Hence, if we feed the model randomly selected windows, it will still treat them as if they were consecutive, and retain the state from one window to the next. To avoid this, we therefore need to make sure that the generator outputs windows of training data that are indeed consecutive (and not randomly selected as above).

2. When training the model we will process the whole training data multiple times (i.e. we train for multiple epochs). However, if we have statefulness between epochs this would effectively result in a "circular dependence", where the final state at time step $t = n_{\text{train}}$ would be used as the initial state at time $t = 1$. To avoid this, we can manually reset the state of the model by calling `model.reset_states()`.

Taking this two points into consideration, complete the code for the stateful data generator below.

```
In [25]: def generator_train_stateful(window_size, model):
    """In addition to the window size, the generator also takes the model as input so
    that we can reset the RNN states at appropriate intervals."""

    # Compute the total number of windows of length window_size that we need to cover all the training data.
    # Note 1: The length of x_train and y_train is ntrain-1 since we work with 1-step prediction.
    # Note 2: The final window should be smaller than window_size, if (ntrain-1) is not evenly divisible by (window_size - 1)
    n_windows = (ntrain-1) // (window_size - 1)

    for i in range(n_windows):
        # First time index of window (inclusive)
        start_of_window = i * window_size
        # Last time index of window (exclusive, i.e. this is the index to the first time step after the window)
        end_of_window = (i+1) * window_size
        yield x_train[start_of_window:end_of_window, :], y_train[start_of_window:end_of_window, :]

    model.reset_states()

    """NOTE: In addition to replacing the ?????? with the correct code, you need to move the line"""
    """to the correct place in the function definition above!"""
```

With the generator defined we can train the model.

```
In [26]: window_size = 100
model3.compile(loss='mse', optimizer='rmsprop', metrics=['mse'])
history = model3.fit(generator_train_stateful(window_size, model3),
                    epochs = 200,
                    batch_size = 1,
                    validation_data = (x_test, y_test),
                    callbacks=[model_checkpoint_callback])

Similarly to above we plot the error curves vs the iteration (epoch) number.
```

```
In [27]: plot_history(history, start_at)
```

Q11: Comparing this error plot to the one you got for training Options 1 and 2, can you see any *qualitative* differences?

Optional: If you have a theory regarding the reason for the observed differences, feel free to explain!

A11: The test error is almost always higher than to the training error even though this wasn't the case in option 2. The reason for this is that with the training data now being continuous the model has prior knowledge and therefore performs better in training.

Q12: Compute a prediction for all values of $\{y_2, \dots, y_n\}$ analogously to **Q6**.

```
In [28]: # Predict on all data using the final model.
y_n = y.reshape((1,ndata-1,1))
# We predict using y_1,...,y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n.
y_pred3 = model3.predict(y_n) # We predict using y_1,...,y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n.
y_pred3 = model3.predict(y_n).flatten() + y[ntrain-1:]

In [29]: # Plot prediction on test data
plot_prediction(y_pred3)

# Evaluate MSE and MAE (both training and test data)
evaluate_performance(y_pred3, y[1:], ntrain-1, name='Stacked RNN, windowing/stateful')

Model Stacked RNN, windowing/stateful
Training MSE: 618.8775, MAE: 17.5515
Testing MSE: 577.8973, MAE: 17.6879

Predictions on test data
```

5. Reflection

Q13: Which model performed best? Did you manage to improve the prediction compared to the two baseline methods? Did the RNN models live up to your expectations? Why/Why not? Please reflect on the lab using a few sentences.

A13: Comparing three models, it looks like Option 3 has least MSE (570.6165). So, it is safe to say Option 3 performed best. RNN models live up to expectations because even when a sequence of data is passed without defining window size or states per epoch, the test MSE was very low. By defining these values we have reduced test MSE from 573.2477 to 570.6165 which is very close. Hence, we can say all the RNN model performed well with the data.

6. A more complex network (OPTIONAL)

If you are interested, feel free to play around with more complex models and see if you can improve the predictive performance! It is very easy to build stacked models in `Keras`, see the example below.

```
In [30]: # A stacked model with 3 layers of LSTM cells, two Dense layers with ReLU activation and a final linear output
model4 = keras.Sequential([
    tf.keras.layers.LSTM(64, batch_input_shape=(1, None, 1), return_sequences=True, stateful=True),
    tf.keras.layers.LSTM(64, batch_input_shape=(1, None, 1), return_sequences=True, stateful=True),
    tf.keras.layers.LSTM(64, batch_input_shape=(1, None, 1), return_sequences=True, stateful=True),
    tf.keras.layers.Dense(32, activation='relu'),
    tf.keras.layers.Dense(16, activation='relu'),
    tf.keras.layers.Dense(1),
])
model4.summary()

Model: "sequential_2"
```

Layer (type)	Output Shape	Param #
lstm (LSTM)	(1, None, 64)	16896
lstm_1 (LSTM)	(1, None, 64)	33024
lstm_2 (LSTM)	(1, None, 64)	33024
dense_2 (Dense)	(1, None, 32)	2080
dense_3 (Dense)	(1, None, 16)	17
dense_4 (Dense)	(1, None, 1)	1
Total params:	85,569	
Trainable params:	85,569	
Non-trainable params:	0	

We can store the best model in a file, so that we can load it after analysing the training procedure.

```
In [31]: checkpoint_filepath = './'
model_checkpoint_callback = tf.keras.callbacks.ModelCheckpoint(
    filepath=checkpoint_filepath,
    save_weights_only=True,
    monitor='val_loss',
    save_best_only=True) # Save only the best model, determined by the validation loss

Train the model
```

```
In [32]: window_size = 100
model4.compile(loss='mse', optimizer='rmsprop', metrics=['mse'])
history = model4.fit(generator_train_stateful(window_size, model4),
                    epochs = 200,
                    batch_size = 1,
                    validation_data = (x_test, y_test),
                    callbacks=[model_checkpoint_callback])

-----
InvalidArgumentError: slice index -1 of dimension 0 out of bounds.
[[node gradients/strided_slice_2_grad/StridedSliceGrad]]
[[node gradients/PartitionedCall]] [[Op::inference_train_function_42849]]
Function call stack:
train_function -> train_function -> train_function
```

Q14 (optional): Based on the training and test error plots, are there signs of over- or underfitting?

A14: We load the best model from checkpoint.

```
In [33]: model4.load_weights(checkpoint_filepath)

In [34]: # Predict on all data using the final model.
y_n = y.reshape((1,ndata-1,1))
# We predict using y_1,...,y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n.
y_pred4 = model4.predict(y_n) # We predict using y_1,...,y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n.
y_pred4 = model4.predict(y_n).flatten() + y[ntrain-1:]

In [35]: # Plot prediction on test data
plot_prediction(y_pred4)

# Evaluate MSE and MAE (both training and test data)
evaluate_performance(y_pred4, y[1:], ntrain-1, name='Stacked RNN, windowing/stateful')

Model Stacked RNN, windowing/stateful
Training MSE: 618.8775, MAE: 17.5515
Testing MSE: 577.8973, MAE: 17.6879

Predictions on test data
```