

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 layers
import matplotlib.pyplot as plt

plt.rcParams["figure.figsize"] = (16,6) # 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 use a temporal forecasting on the data of [sunspots](#). We work with a data set that has been published on [Kaggle](#), with the description:

Sunspots are temporary 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 2017-06-31.

```
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
ntrain=nlen(y)
print(f'Total number of data points: {ntrain}')

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

Total number of data points: 325276
Number of training data points: 227692

In [4]: plt.plot(dates[ntrain:], y[ntrain:])
plt.plot(dates[ntrain:], y[ntrain:])
plt.xticks(range(0, n_test, 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 very 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 the constant into account in the lab assignment and allow ourselves to model the data using a Gaussian likelihood (i.e. using MSE as 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: 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 errors in prediction
    resid = y - y_pred

    # 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(abs(resid[split_time:]))*MAX_VAL
    test_mae = np.mean(abs(resid[split_time:]))*MAX_VAL

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

Q1: Implement the naive baseline method which predicts according to  $\hat{y}_{t+1} = y_t$ . Since the previous value is needed for the prediction we do not get a prediction at  $t = 1$ . Hence, we evaluate the method by predicting values at  $t = 2, \dots, (n)$  (cf. an AR(p) model where we start predicting at  $t = p + 1$ ).
```

```
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[i] is a prediction of y[i+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 important that
# you understand it!
y_pred_naive = y[:-1]

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

Model naive
Training MSE: 776.5437, MAE: 19.3285
Testing MSE: 788.6369, MAE: 19.2256
Note 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

p=30 # 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[ntrain:], ..., y[n]-1.
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-1, # Number of leading elements in the input arrays corresponding to training data points
                    name='AR')

Model AR
Training MSE: 683.8856, MAE: 17.3426
Testing MSE: 590.2732, MAE: 17.0221
```

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 $[Q, H, P]$ 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 terms 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,M,P] = [batch, timesteps, feature] = [1, 3252, 1]`

- Q : The number of sequences
- M = n_{data} : The length of the sequences (number of time steps)
- P : The number of features (at each time step), i.e. the dimension of each time series

Q3: Continuing the question above, answer the following:

- What is the meaning of setting `units = d`?
- Assume 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 hard set the parameter `return_sequences=False` when constructing the layer, then what would be the answer to the previous question?

- A3:
- d is the number of hidden layers that pass information from one hidden state to the next i hidden state dimension
 - $[1,n,d]$
 - $[1,d]$

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,

$$\mathbf{h}_t = \sigma(\mathbf{W}\mathbf{h}_{t-1} + \mathbf{U}\mathbf{g}_{t-1} + \mathbf{b}),$$
$$\hat{\mathbf{y}}_{t+1} = \mathbf{C}\mathbf{h}_t + \mathbf{c},$$

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

model0=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
model0.summary()
init_weights = model0.get_weights().copy()

Model: "sequential"

Layer (type) Output Shape Param #
=====
simple_rnn (SimpleRNN) (None, None, 10) 120
Dense (Dense) (None, None, 1) 11
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:
•  $W = d \times d = 10 \times 10 = 100$ 
•  $U = d \times \text{input\_shape} = 10 \times 1 = 10$ 
•  $b = d \times 1 = 10 \times 1 = 10$ 
•  $W = U + b = 120$  which is the number of parameters of the RNN cells
•  $C = d \times 1$  (dense layer unit) = 10
•  $c = 1$  (dense layer unit)
•  $C + c = 11$  which is the number of parameters of the prediction dense layer
• or param\_number = output\_channel\_number (input\_channel\_number + 1) = 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 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 \hat{y}_{t+1} . Equivalently, we can define the target value at time step t to be the residual $y_t := y_t - y_{t-1}$. 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[t][0]=y[t-1], ..., x[ntrain-1]=y[ntrain-1]
y_train = y[ntrain:] # Output is denoted by yt, training outputs are y[t][0]=y[t]-y[t-1], ..., y[ntrain-1]=y[ntrain]-y[ntrain-1]

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

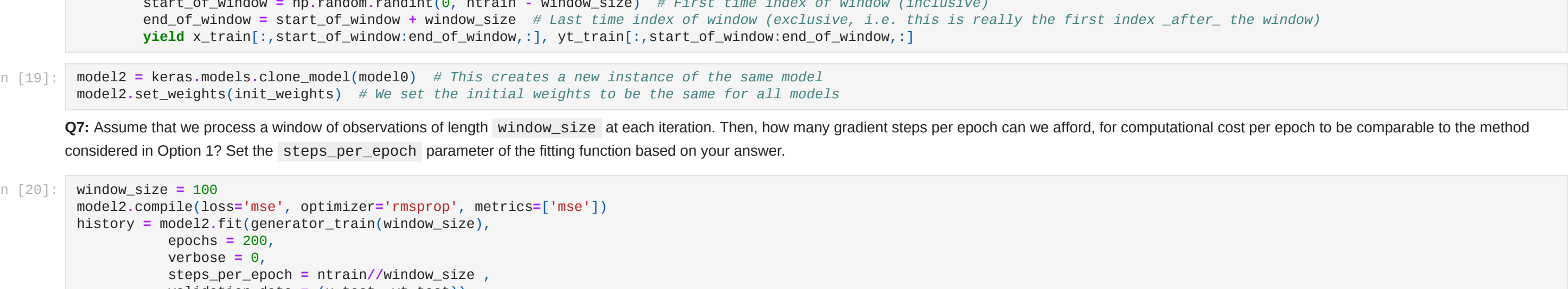
```
In [12]: model1 = keras.models.clone_model(model0) # This creates a new instance of the same model
model1.set_weights(init_weights) # We set the initial weights to be the same for all models

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 = len(y_train),
                    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.
```



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 as $\{y_t\}$ to the `predict` to test data using the input shape used in Keras (cf. above).

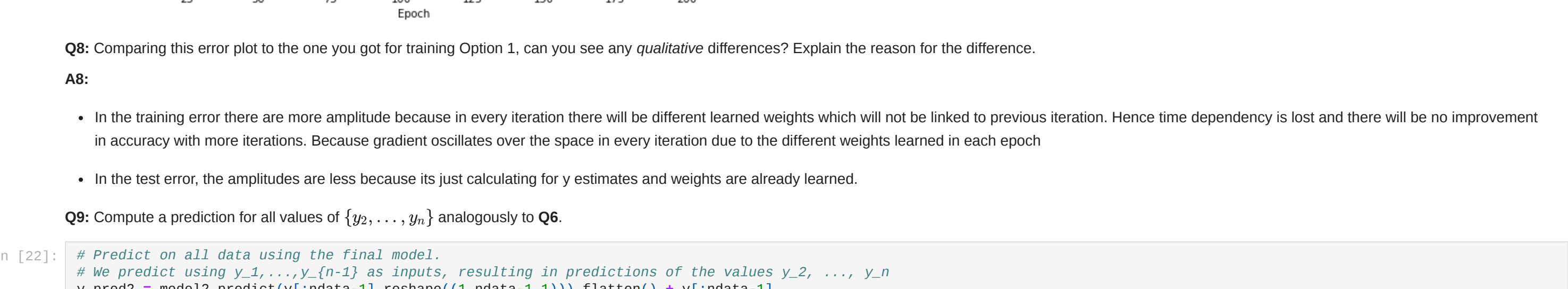
Hint: Since the model is trained on the residuals y_t , don't forget to add back y_{t-1} when predicting y_t . However, make sure that you don't "cheat" by using a non-causal predictor (i.e. using y_t when predicting y_t)!

```
In [15]: # Predict on all data using the final model.

# We predict using y_1, ..., y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n.
# That is, y_pred should be an (n-1) array where element y_pred[i] is based only on values y[1:]
# Note: The final window could be smaller than window_size, if (ntrain-1) is not evenly divisible by the window_size.
number_of_windows = int(ntrain-1/window_size)

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_predictions(y_pred):
    # Plot prediction on test data
    plt.plot(dates[ntrain:], y[ntrain:])
    plt.plot(dates[ntrain:], y_pred[ntrain-1:])
    plt.xticks(range(0, n_test, 300), rotation = 90); # Show only one tick every 25th year for clarity
    plt.legend(['Data', 'Prediction'])
    plt.title('Predictions on test data')
```



Q7: 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.

- In the training error there are more amplitude because in every iteration there is a different learned weights which will not be learned to previous. Hence time dependency is lost and there will be no improvement in accuracy with more iterations. Because gradient oscillates over the space in every iteration due to the different weights and biases are already learned.
- In the test error, the amplitudes are less because it is just calculating for y estimates and weights are already learned.

Q8: Compute a prediction for all values of $\{y_1, \dots, y_n\}$ analogously to Q6.

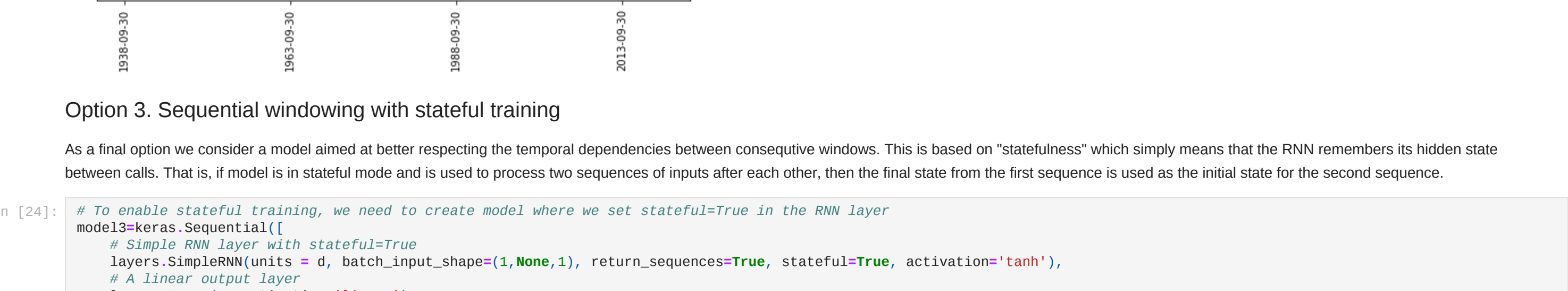
```
In [22]: # Predict on all data using the final model.
# We predict using y_1, ..., y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n
y_pred2 = model2.predict(y[ntrain-1:].reshape((1,ntrain-1,1)).flatten()) + y[ntrain-1:]

In [23]: # Plot prediction on test data
plt.plot(dates[ntrain:], y[ntrain:])
plt.plot(dates[ntrain:], y_pred2[ntrain-1:])

# 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: 602.6192, MAE: 17.5562
Testing MSE: 570.8144, MAE: 17.3017

Predictions on test data
```



Option 2. Random windowing

Instead of using all the training data and 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:
        # The upper value is excluded in randint, so the maximum value that we can get is t = ntrain-window_size-1.
        # Note: The length of x_train (and y_train) is ntrain-1, in agreement with the fact that the size of input/output is ntrain-1.
        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) - this is really the first index after the window
        yield x_train[start_of_window:end_of_window-1], y_train[start_of_window:end_of_window-1]

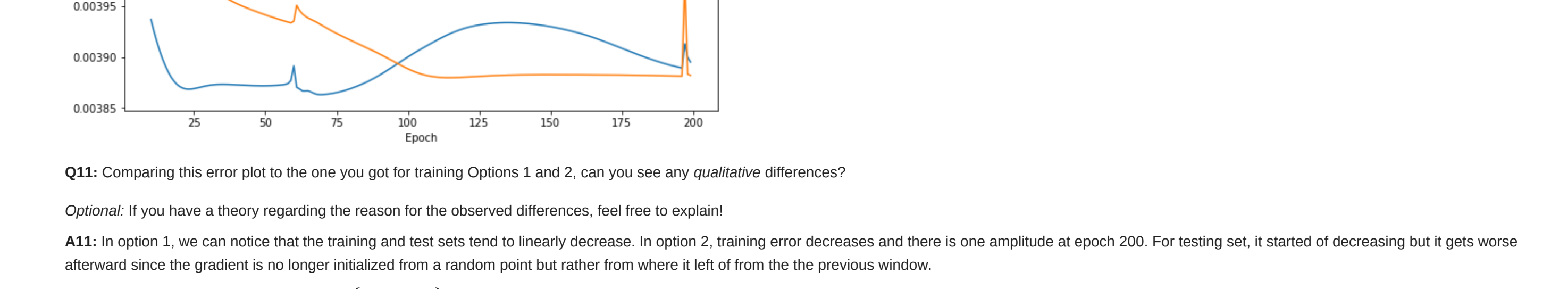
In [19]: model2 = keras.models.clone_model(model0) # This creates a new instance of the same model
model2.set_weights(init_weights) # We set the initial weights to be the same for all models

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 based on your answer.
```

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

Training MSE: 624.8471, MAE: 17.8758
Testing MSE: 616.7663, MAE: 18.8701
```

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



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: In option 1, we can note that the training and test sets tend to linearly decrease. In option 2, training error decreases and there is one amplitude at epoch 200. For testing set, it started of decreasing but it gets worse afterward since the gradient is no longer initialized from a random point but rather from where it left from the previous window and is analogous to Q6.

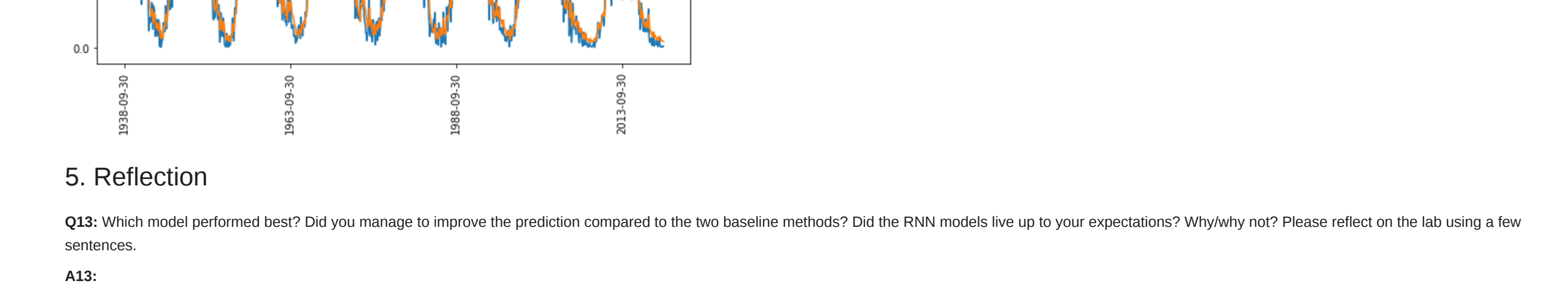
```
In [22]: # Predict on all data using the final model.
# We predict using y_1, ..., y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n
y_pred3 = model3.predict(y[ntrain-1:].reshape((1,ntrain-1,1)).flatten()) + y[ntrain-1:]

In [23]: # Plot prediction on test data
plt.plot(dates[ntrain:], y[ntrain:])
plt.plot(dates[ntrain:], y_pred3[ntrain-1:])

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

Model Simple RNN, windowing/stateful
Training MSE: 624.8471, MAE: 17.8758
Testing MSE: 616.7663, MAE: 18.8701

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 consecutive windows. This is based on "statefulness" in the RNN, meaning that the RNN remembers its hidden state between calls. That is, if the 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(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 of the model, we need to reset the state of the model by calling `model.reset_states()`, where the final state at time step $t = n_{train}$ 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: The length of x_train (and y_train) is ntrain-1, so we work with ntrain-1.
    number_of_windows = int(ntrain-1/window_size)

    while True:
        for i in range(number_of_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 = start_of_window + window_size

            # Note: Python allows using end_of_window = ntrain-1, it will simply truncate the indexing at the final element of the array!

            yield x_train[start_of_window:end_of_window-1], y_train[start_of_window:end_of_window-1]

            model.reset_states()

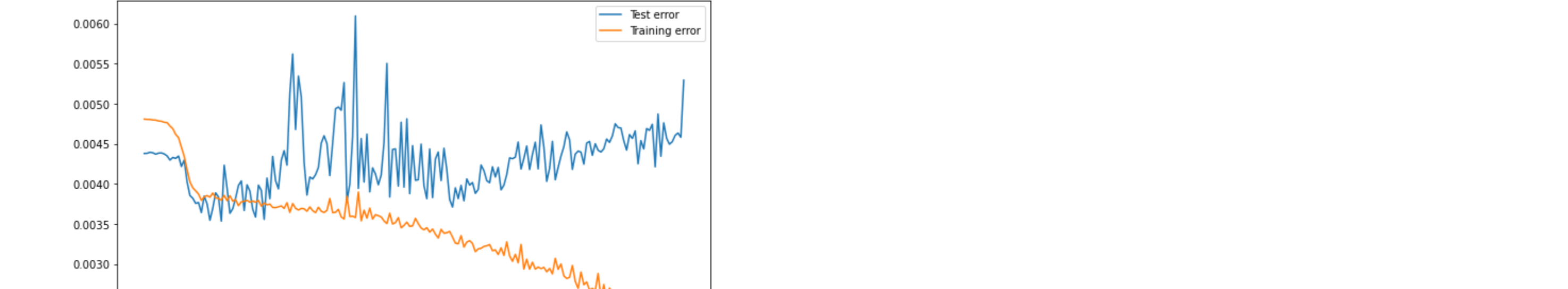
            # NOTE: In addition to replacing the ???? with the correct code, you need to move the line!""
            """to the correct place in the function definition above!"""
```

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,
                    verbose = 0,
                    steps_per_epoch = ntrain/window_size,
                    validation_data = (x_test, y_test))

Training MSE: 624.8471, MAE: 17.8758
Testing MSE: 616.7663, MAE: 18.8701
```

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



Q12: If you have a theory regarding the reason for the observed differences, feel free to explain!

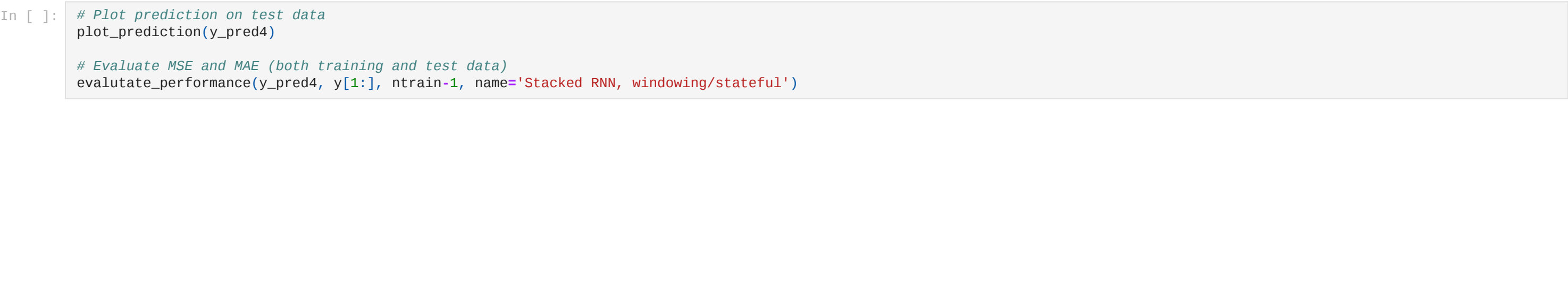
A12: In option 1, we can note that the training and test sets tend to linearly decrease. In option 2, training error decreases and there is one amplitude at epoch 200. For testing set, it started of decreasing but it gets worse afterward since the gradient is no longer initialized from a random point but rather from where it left from the previous window and is analogous to Q6.

```
In [22]: # Predict on all data using the final model.
# 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[ntrain-1:].reshape((1,ntrain-1,1)).flatten()) + y[ntrain-1:]

In [23]: # Plot prediction on test data
plt.plot(dates[ntrain:], y[ntrain:])
plt.plot(dates[ntrain:], y_pred4[ntrain-1:])

# 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: 624.8471, MAE: 17.8758
Testing MSE: 616.7663, MAE: 18.8701
```



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 [35]: model4.load_weights(checkpoint_filepath)

<tensorflow.python.training.tracking.util.CheckpointLoadStatus at 8x258a8325e8>

In [36]: # Predict on all data using the final model.
# We predict using y_1, ..., y_{n-1} as inputs, resulting in predictions of the values y_2, ..., y_n
y_pred4 = model4.predict(?????) # ????? = ?????

File "<ipython-input-36-3ebcd0f32eca>", line 1
      y_pred4 = model4.predict(?????) # ????? = ?????
SyntaxError: invalid syntax

In [37]: # Predict on all data using the final model.
# 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[ntrain-1:].reshape((1,ntrain-1,1)).flatten()) + y[ntrain-1:]

In [38]: # Plot prediction on test data
plt.plot(dates[ntrain:], y[ntrain:])
plt.plot(dates[ntrain:], y_pred4[ntrain-1:])

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