

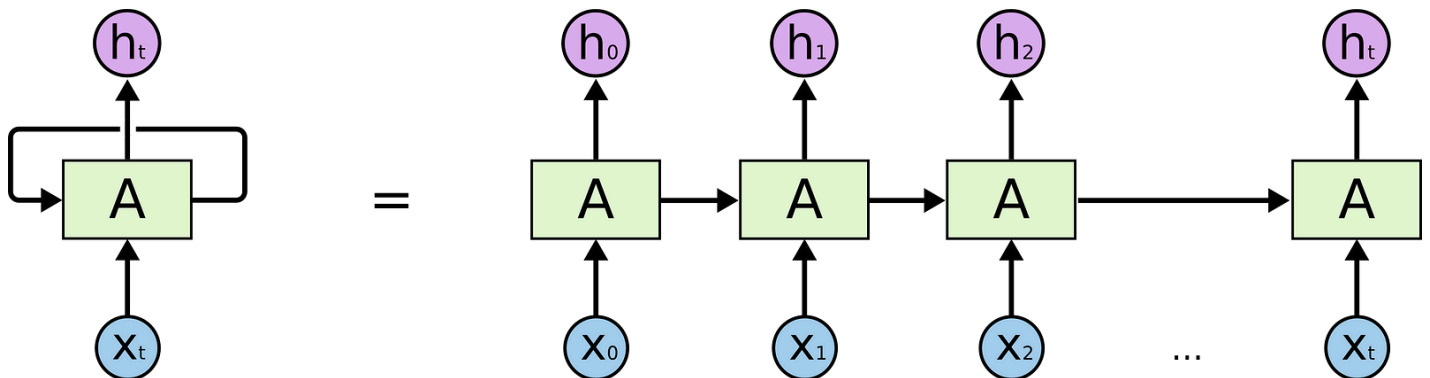
✓ LAB 02 RNN - LSTM

Aim and motivation

The primary reason we have chosen to create this kernel is to practice and use RNNs for various tasks and applications. First of which is time series data. RNNs have truly changed the way sequential data is forecasted. My goal here is to create the ultimate reference for RNNs here on kaggle.

✓ Recurrent Neural Networks

In a recurrent neural network we store the output activations from one or more of the layers of the network. Often these are hidden layer activations. Then, the next time we feed an input example to the network, we include the previously-stored outputs as additional inputs. You can think of the additional inputs as being concatenated to the end of the "normal" inputs to the previous layer. For example, if a hidden layer had 10 regular input nodes and 128 hidden nodes in the layer, then it would actually have 138 total inputs (assuming you are feeding the layer's outputs into itself à la Elman) rather than into another layer). Of course, the very first time you try to compute the output of the network you'll need to fill in those extra 128 inputs with 0s or something.



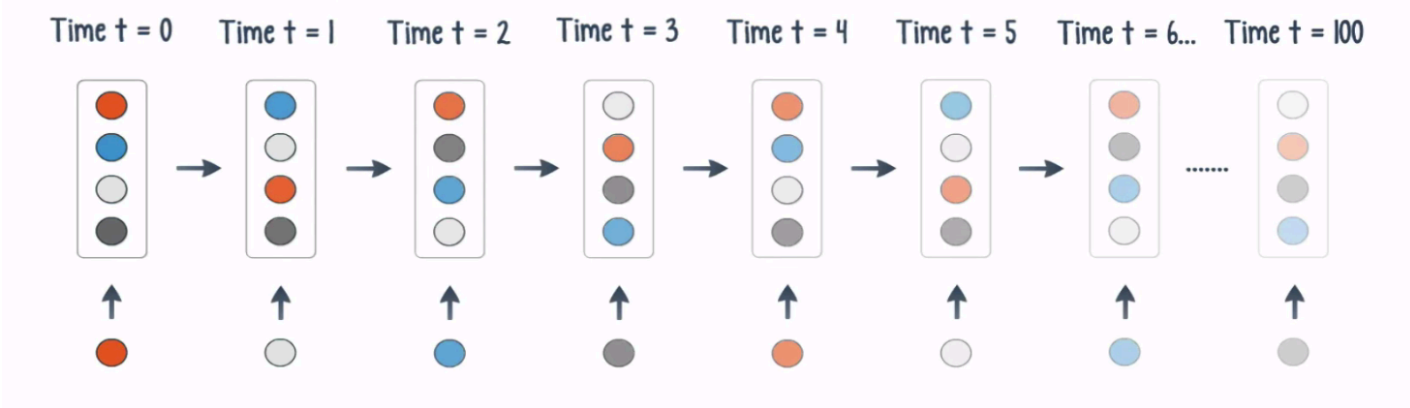
Now, even though RNNs are quite powerful, they suffer from **Vanishing gradient problem** ** which hinders them from using long term information, like they are good for storing memory 3-4 instances of past iterations but larger number of instances don't provide good results so we don't just use regular RNNs. Instead, we use a better variation of RNNs: **Long Short Term Networks(LSTM).

What is Vanishing Gradient problem?

Vanishing gradient problem is a difficulty found in training artificial neural networks with gradient-based learning methods and backpropagation. In such methods, each of the neural network's weights receives an update proportional to the partial derivative of the error function with respect to the current weight in each iteration of training. The problem is that in some cases, the gradient will be vanishingly small, effectively preventing the weight from changing its value. In the worst case, this may completely stop the neural network from further training. As one example of the problem cause, traditional activation functions such as the hyperbolic tangent function have gradients in the range (0, 1), and backpropagation computes gradients by the chain rule. This has the effect of multiplying n of these small numbers to compute gradients of the "front" layers in an n-layer network, meaning that the gradient (error signal) decreases exponentially with n while the front layers train very slowly.

Source: [Wikipedia](https://en.wikipedia.org/wiki/Vanishing_gradient_problem)

Decay of information through time

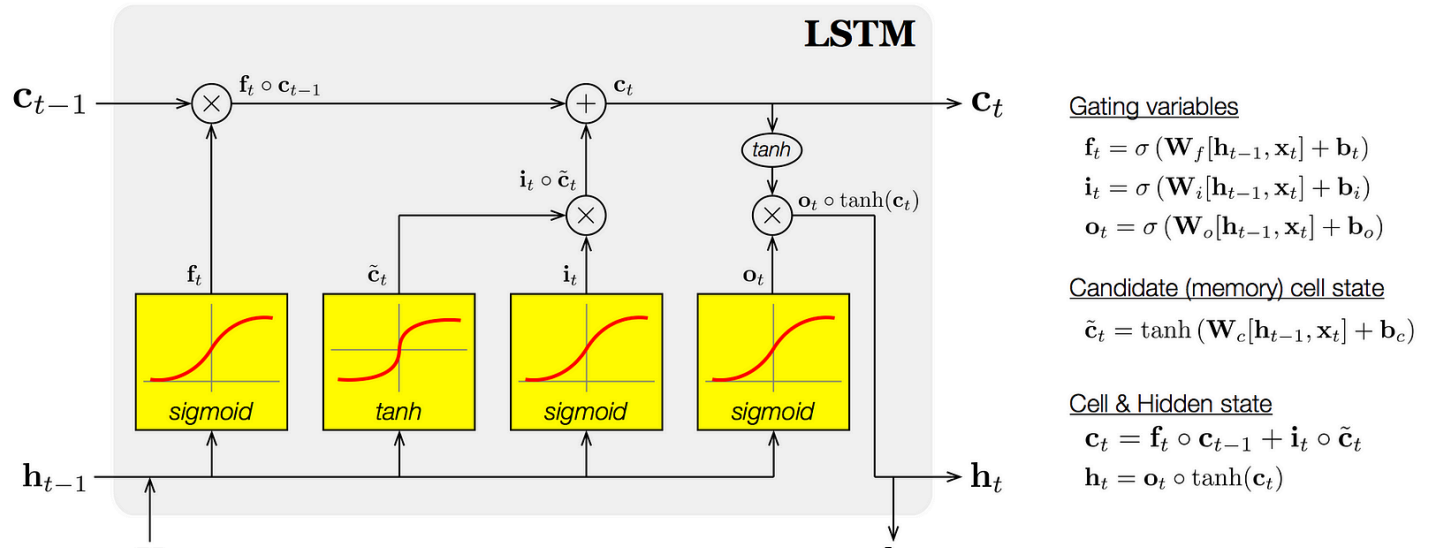


Long Short Term Memory(LSTM)

Long short-term memory (LSTM) units (or blocks) are a building unit for layers of a recurrent neural network (RNN). A RNN composed of LSTM units is often called an LSTM network. A common LSTM unit is composed of a cell, an input gate, an output gate and a forget gate. The cell is responsible for "remembering" values over arbitrary time intervals; hence the word "memory" in LSTM. Each of the three gates can be thought of as a "conventional" artificial neuron, as in a multi-layer (or feedforward) neural network: that is, they compute an activation (using an activation function) of a weighted sum. Intuitively, they can be thought as regulators of the flow of values that goes through the connections of the LSTM; hence the denotation "gate". There are connections between these gates and the cell.

The expression long short-term refers to the fact that LSTM is a model for the short-term memory which can last for a long period of time. An LSTM is well-suited to classify, process and predict time series given time lags of unknown size and duration between important events. LSTMs were developed to deal with the exploding and vanishing gradient problem when training traditional RNNs.

Source: [Wikipedia](#)



Components of LSTMs

So the LSTM cell contains the following components

- Forget Gate "f" (a neural network with sigmoid)
- Candidate layer "C"(a NN with Tanh)
- Input Gate "I" (a NN with sigmoid)
- Output Gate "O"(a NN with sigmoid)
- Hidden state "H" (a vector)
- Memory state "C" (a vector)
- Inputs to the LSTM cell at any step are X_t (current input) , H_{t-1} (previous hidden state) and C_{t-1} (previous memory state).
- Outputs from the LSTM cell are H_t (current hidden state) and C_t (current memory state)

✓ Working of gates in LSTMs

First, LSTM cell takes the previous memory state C_{t-1} and does element wise multiplication with forget gate (f) to decide if present memory state C_t . If forget gate value is 0 then previous memory state is completely forgotten else if forget gate value is 1 then previous memory state is completely passed to the cell (Remember f gate gives values between 0 and 1).

$$C_t = C_{t-1} * f_t$$

Calculating the new memory state:

$$C_t = C_t + (i_t * C'_t)$$

Now, we calculate the output:

$$H_t = \tanh(C_t)$$

✓ And now we get to the code...

I will use LSTMs for predicting the price of stocks of IBM for the year 2017

```
1 # Importing the libraries
2 import numpy as np
3 import matplotlib.pyplot as plt
4 plt.style.use('fivethirtyeight')
5 import pandas as pd
6 from sklearn.preprocessing import MinMaxScaler
7 from keras.models import Sequential
8 from keras.layers import Dense, LSTM, Dropout, GRU, Bidirectional
9 from keras.optimizers import SGD
10 import math
11 from sklearn.metrics import mean_squared_error

1 # Some functions to help out with
2 def plot_predictions(test,predicted):
3     plt.plot(test, color='red',label='Real IBM Stock Price')
4     plt.plot(predicted, color='blue',label='Predicted IBM Stock Price')
5     plt.title('IBM Stock Price Prediction')
6     plt.xlabel('Time')
7     plt.ylabel('IBM Stock Price')
8     plt.legend()
9     plt.show()
10
11 def return_rmse(test,predicted):
12     rmse = math.sqrt(mean_squared_error(test, predicted))
13     print("The root mean squared error is {}".format(rmse))

1 # First, we get the data
2 dataset = pd.read_csv('../input/IBM_2006-01-01_to_2018-01-01.csv', index_col='Date', pars
```

```
3 dataset.head()
```

```
1 # Checking for missing values
2 training_set = dataset[:, '2016'].iloc[:, 1:2].values
3 test_set = dataset[:, '2017:'].iloc[:, 1:2].values
```

```
1 # We have chosen 'High' attribute for prices. Let's see what it looks like
2 dataset["High"][:, '2016'].plot(figsize=(16,4), legend=True)
3 dataset["High"][:, '2017:'].plot(figsize=(16,4), legend=True)
4 plt.legend(['Training set (Before 2017)', 'Test set (2017 and beyond)'])
5 plt.title('IBM stock price')
6 plt.show()
```

```
1 # Scaling the training set
2 sc = MinMaxScaler(feature_range=(0,1))
3 training_set_scaled = sc.fit_transform(training_set)
```

```
1 # Since LSTMs store long term memory state, we create a data structure with 60 timesteps
2 # So for each element of training set, we have 60 previous training set elements
3 X_train = []
4 y_train = []
5 for i in range(60, 2769):
6     X_train.append(training_set_scaled[i-60:i, 0])
7     y_train.append(training_set_scaled[i, 0])
8 X_train, y_train = np.array(X_train), np.array(y_train)
```

```
1 # Reshaping X_train for efficient modelling
2 X_train = np.reshape(X_train, (X_train.shape[0], X_train.shape[1], 1))
```

```
1 # The LSTM architecture
2 regressor = Sequential()
3 # First LSTM layer with Dropout regularisation
4 regressor.add(LSTM(units=50, return_sequences=True, input_shape=(X_train.shape[1], 1)))
5 regressor.add(Dropout(0.2))
6 # Second LSTM layer
7 regressor.add(LSTM(units=50, return_sequences=True))
8 regressor.add(Dropout(0.2))
9 # Third LSTM layer
10 regressor.add(LSTM(units=50, return_sequences=True))
11 regressor.add(Dropout(0.2))
12 # Fourth LSTM layer
13 regressor.add(LSTM(units=50))
14 regressor.add(Dropout(0.2))
15 # The output layer
16 regressor.add(Dense(units=1))
17
18 # Compiling the RNN
```

```

19 regressor.compile(optimizer='rmsprop',loss='mean_squared_error')
20 # Fitting to the training set
21 regressor.fit(X_train,y_train,epochs=50,batch_size=32)

1 # Now to get the test set ready in a similar way as the training set.
2 # The following has been done so forst 60 entires of test set have 60 previous values whi
3 # 'High' attribute data for processing
4 dataset_total = pd.concat((dataset["High"][:'2016'],dataset["High"][:'2017':]),axis=0)
5 inputs = dataset_total[len(dataset_total)-len(test_set) - 60:].values
6 inputs = inputs.reshape(-1,1)
7 inputs = sc.transform(inputs)

1 # Preparing X_test and predicting the prices
2 X_test = []
3 for i in range(60,311):
4     X_test.append(inputs[i-60:i,0])
5 X_test = np.array(X_test)
6 X_test = np.reshape(X_test, (X_test.shape[0],X_test.shape[1],1))
7 predicted_stock_price = regressor.predict(X_test)
8 predicted_stock_price = sc.inverse_transform(predicted_stock_price)

1 # Visualizing the results for LSTM
2 plot_predictions(test_set,predicted_stock_price)

1 # Evaluating our model
2 return_rmse(test_set,predicted_stock_price)

```

Truth be told. That's one awesome score.

LSTM is not the only kind of unit that has taken the world of Deep Learning by a storm. We have **Gated Recurrent Units(GRU)**. It's not known, which is better: GRU or LSTM because they have comparable performances. GRUs are easier to train than LSTMs.

✓ Gated Recurrent Units

In simple words, the GRU unit does not have to use a memory unit to control the flow of information like the LSTM unit. It can directly makes use of the all hidden states without any control. GRUs have fewer parameters and thus may train a bit faster or need less data to generalize. But, with large data, the LSTMs with higher expressiveness may lead to better results.

They are almost similar to LSTMs except that they have two gates: reset gate and update gate. Reset gate determines how to combine new input to previous memory and update gate determines how much of the previous state to keep. Update gate in GRU is what input gate and forget gate were

in LSTM. We don't have the second non linearity in GRU before calculating the output, neither they have the output gate.



```
1 # The GRU architecture
2 regressorGRU = Sequential()
3 # First GRU layer with Dropout regularisation
4 regressorGRU.add(GRU(units=50, return_sequences=True, input_shape=(X_train.shape[1],1), a
5 regressorGRU.add(Dropout(0.2))
6 # Second GRU layer
7 regressorGRU.add(GRU(units=50, return_sequences=True, input_shape=(X_train.shape[1],1), a
8 regressorGRU.add(Dropout(0.2))
9 # Third GRU layer
10 regressorGRU.add(GRU(units=50, return_sequences=True, input_shape=(X_train.shape[1],1), a
11 regressorGRU.add(Dropout(0.2))
12 # Fourth GRU layer
13 regressorGRU.add(GRU(units=50, activation='tanh'))
14 regressorGRU.add(Dropout(0.2))
15 # The output layer
16 regressorGRU.add(Dense(units=1))
17 # Compiling the RNN
18 regressorGRU.compile(optimizer=SGD(lr=0.01, decay=1e-7, momentum=0.9, nesterov=False), los
19 # Fitting to the training set
20 regressorGRU.fit(X_train,y_train,epochs=50,batch_size=150)
```

The current version uses a dense GRU network with 100 units as opposed to the GRU network with 50 units in previous version

```
1 # Preparing X_test and predicting the prices
2 X_test = []
3 for i in range(60,311):
4     X_test.append(inputs[i-60:i,0])
5 X_test = np.array(X_test)
6 X_test = np.reshape(X_test, (X_test.shape[0],X_test.shape[1],1))
7 GRU_predicted_stock_price = regressorGRU.predict(X_test)
8 GRU_predicted_stock_price = sc.inverse_transform(GRU_predicted_stock_price)
```

```
1 # Visualizing the results for GRU
2 plot_predictions(test_set,GRU_predicted_stock_price)
```

```
1 # Evaluating GRU
2 return_rmse(test_set,GRU_predicted_stock_price)
```

✓ Sequence Generation

Here, I will generate a sequence using just initial 60 values instead of using last 60 values for every new prediction. **Due to doubts in various comments about predictions making use of test set values, I have decided to include sequence generation.** The above models make use of test set so it is using last 60 true values for predicting the new value(I will call it a benchmark). This is why the error is so low. Strong models can bring similar results like above models for sequences too but they require more than just data which has previous values. In case of stocks, we need to know the sentiments of the market, the movement of other stocks and a lot more. So, don't expect a remotely accurate plot. The error will be great and the best I can do is generate the trend similar to the test set.

I will use GRU model for predictions. You can try this using LSTMs also. I have modified GRU model above to get the best sequence possible. I have run the model four times and two times I got error of around 8 to 9. The worst case had an error of around 11. Let's see what this iterations.

The GRU model in the previous versions is fine too. Just a little tweaking was required to get good sequences. **The main goal of this kernel is to show how to build RNN models. How you predict data and what kind of data you predict is up to you. I can't give you some 100 lines of code where you put the destination of training and test set and get world-class results. That's something you have to do yourself.**

```
1 # Preparing sequence data
2 initial_sequence = X_train[2708,:]
3 sequence = []
4 for i in range(251):
5     new_prediction = regressorGRU.predict(initial_sequence.reshape(initial_sequence.shape
6     initial_sequence = initial_sequence[1:]
7     initial_sequence = np.append(initial_sequence,new_prediction,axis=0)
8     sequence.append(new_prediction)
9 sequence = sc.inverse_transform(np.array(sequence).reshape(251,1))
```

```
1 # Visualizing the sequence
2 plot_predictions(test_set,sequence)
```

```
1 # Evaluating the sequence
2 return_rmse(test_set,sequence)
```

So, GRU works better than LSTM in this case. Bidirectional LSTM is also a good way so make the model stronger. But this may vary for different data sets. **Applying both LSTM and GRU together gave even better results.**

