





- ▶ Stock price prediction
 - ▶ A correct stock price prediction requires a lot of past information; it's impossible to predict the price evolution without knowing its previous trends

- ▶ Text generation
 - ▶ Text classification
 - **▶** Translation
 - ▶ Chatbot
 - **)** ...



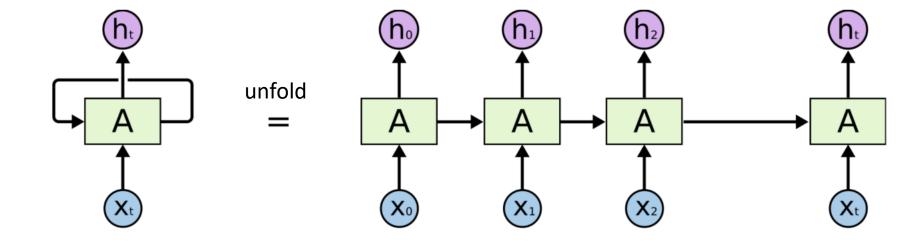
- ▶ Feedforward networks don't consider temporal states
 - Sequences must be passed entirely, not word by word
- A RNN processes sequences by iterating over the elements and by maintaining a state containing information about what has been seen so far

2023/2024 MACHINE LEARNING



Recurrent Neural Networks

- A RNN is a graph with cycles
 - ▶ Perceptrons can benefit from feedback loops
 - \blacktriangleright The output of a perceptron at time t is involved in its input at time t+1





Elman network

$$h_t = \sigma_h(W_h x_t + U_h h_{t-1} + b_h)$$

$$y_t = \sigma_y(W_y h_t + b_y)$$

Jordan network

$$h_t = \sigma_h(W_h x_t + U_h y_{t-1} + b_h)$$
$$y_t = \sigma_y(W_y h_t + b_y)$$

 x_t : input vector

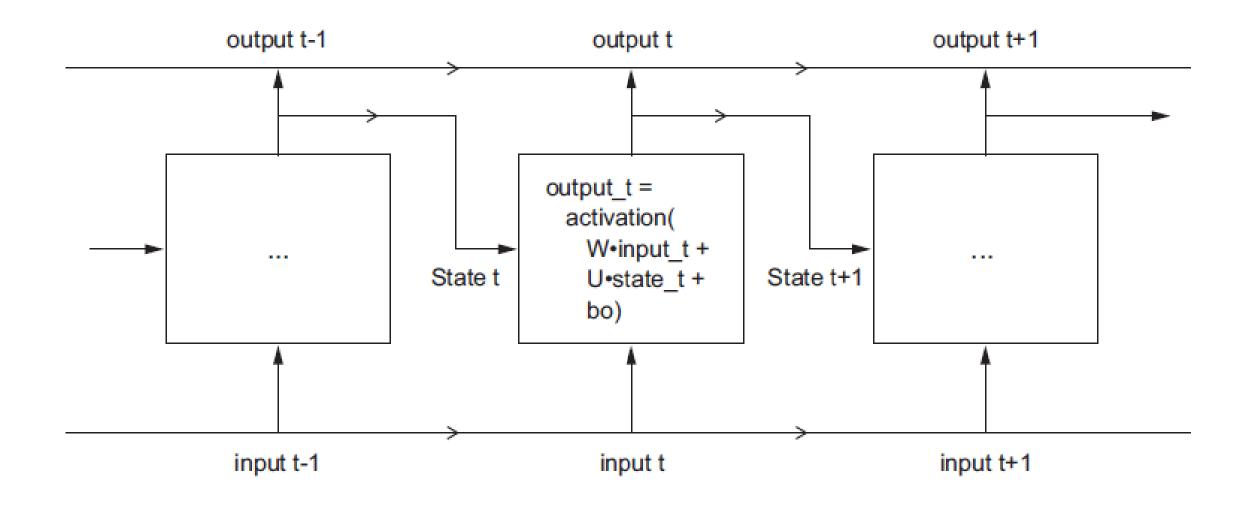
 y_t : output vector

 h_t : hidden layer vector

W, U, b: weights and bias

 σ_h , σ_y : activation functions

Unroll RNN





Implementation

```
Number of timesteps in
                                       Dimensionality of the
 the input sequence
                                       input feature space
     import numpy as np
                                                                           Input data: random
                                                                           noise for the sake of
                                          Dimensionality of the
     timesteps = 100
                                                                           the example
                                          output feature space
     input_features = 32
     output_features = 64
                                                                              Initial state: an
     inputs = np.random.random((timesteps, input_features)) <
                                                                              all-zero vector
     state_t = np.zeros((output_features,))
    W = np.random.random((output_features, input_features))
                                                                             Creates random
    U = np.random.random((output_features, output_features))
                                                                             weight matrices
    b = np.random.random((output_features,))
                                                       input t is a vector of
     successive outputs = []
                                                       shape (input features,).
     for input_t in inputs:
         output_t = np.tanh(np.dot(W, input_t) + np.dot(U, state_t) + b)
         successive_outputs.append(output_t)
         state_t = output_t
     final_output_sequence = np.concatenate(successive_outputs, axis=0) <-
                                                            The final output is a 2D tensor of
 Stores this output in a list
                                                         shape (timesteps, output features).
Combines the input with the current
state (the previous output) to obtain
                                                                       Updates the state of the
                                                                  network for the next timestep
the current output
```



```
model.add(SimpleRNN(32, return_sequences=True))
model.add(SimpleRNN(32, return_sequences=True))
model.add(SimpleRNN(32, return_sequences=True))
model.add(SimpleRNN(32)) # This last layer only returns the last output
model.summary()
```

Example

```
from keras.datasets import imdb
from keras.preprocessing import sequence

max_features = 10000  # number of words to consider as features
maxlen = 500
batch_size = 32

(input_train, y_train), (input_test, y_test) = imdb.load_data(num_words=max_features)
input_train = sequence.pad_sequences(input_train, maxlen=maxlen)
input_test = sequence.pad_sequences(input_test, maxlen=maxlen)
```

Example



Backpropagation Through Time

- ▶ Similarly to any neural network, RNNs needs to learn the weight values according to a training set
- The learning algorithm, namely backpropagation, needs a slightly change due to the feedback loops
- ▶ Its variant is called Backpropagation through time (BPTT)



Backpropagation Through Time

- The idea is very simple:
 - Unfold the recurrent neural network in time
 - ▶ The unfolded neurons share the same parameters
 - Backpropagate the loss gradient like in a feedforward neural network

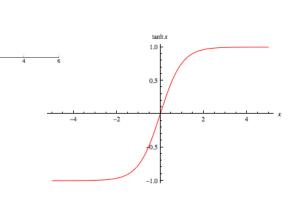
■ Vanishing gradient

According to the number of time steps, an unfolded RNN may be a very long sequence

▶ To simulate neuron activations, perceptrons are typically equipped with close-to-step activation functions:

$$\blacktriangleright \text{ Sigmoid, } \sigma(x) = \frac{1}{1 + e^{-x}}$$

▶ Hyperbolic tangent, $tanh(x) = \frac{e^{x^{2}} - e^{x^{2}}}{e^{x} + e^{-x}}$



Vanishing gradient

▶ Stochastic gradient descent:

$$w^* = w - \eta \nabla E_i(o, f(i, w))$$

- where:
 - $\blacktriangleright w^*$, the new weights
 - ▶ w, the old weights
 - \blacktriangleright η , the learning rate
 - \blacktriangleright $E_i(\cdot)$ the error (loss) function of the j-th observation
 - *f* , the activation function
 - $\blacktriangleright \nabla$, the gradient operator
 - ▶ *i*, the input
 - *o*, the output

Vanishing gradient

In backpropagation we have:

$$\nabla E_j(o, f_l(i, w_l)) = \nabla E_j(o, f_l(f_{l-1}(f_{l-2}(f_{l-3}(...), w_{l-2}), w_{l-1}), w_l))$$

where l is a layer out of L

▶ For a recurrent perceptron (feedback loop)

$$f_l = f_{l-1} = f_{l-2} = \dots = f$$

■ Vanishing gradient

▶ Hence a RNN perceptron has gradient:

$$\nabla E_j(o, f(f(f(\dots), w), w), w))$$

▶ By chain rule of derivatives we know:

$$\frac{d}{dx}f(g(x)) = f'(g(x)) \cdot g'(x)$$

▶ This means that in a RNN we need to multiply the same derivative function many times:

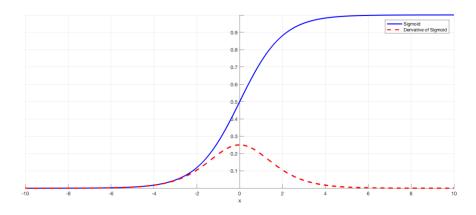
$$[f'(i)]^T$$

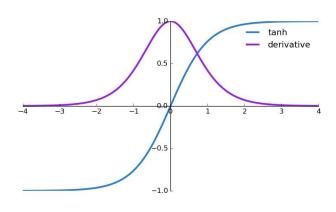
▶ where *T* is the number of time steps

■ Vanishing gradient

▶ The derivative of the sigmoid is:

The derivative of the tanh is:





- In both cases we have quantities that are at most 1
- In $[f'(i)]^T$, let's assume that f'(i) = 0.99 and T = 300

$$0.99^{300} = 0.049$$

- ▶ The gradient is close to 0, it is disappearing... it is vanishing!
 - ▶ The network is not able to find out the optimal values for the weights

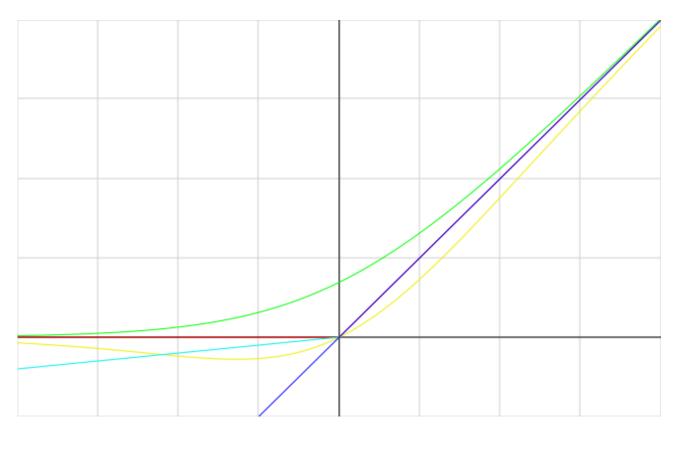


Possible solutions

- ▶ Changing the activation function
 - Linear
 - \blacktriangleright Linear(x) = x
 - ▶ REctified Linear Unit (RELU)

$$RELU(x) = \begin{cases} x & if \ x \ge 0 \\ 0 & otherwise \end{cases}$$

- ▶ Leaky ReLU
 - $LRELU(x; \alpha \ll 1) = \begin{cases} x & if \ x \geq 0 \\ \alpha x & otherwise \end{cases}$
- **▶** Softplus
 - \blacktriangleright Softplus(x) = ln(1 + e^x)
- Swish
 - $\blacktriangleright \text{ Swish}(x) = x \cdot \sigma(x)$
- ▶ There is no control in what to remember
 - ▶ The long-term dependency problem



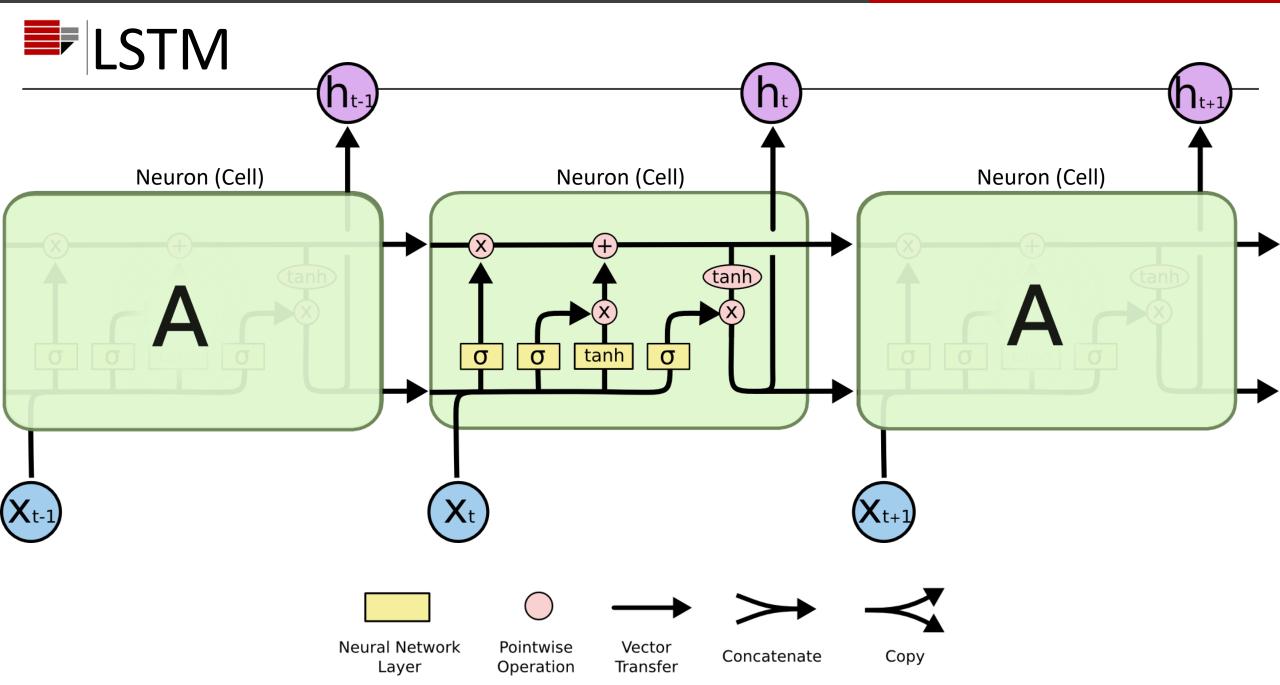
2023/2024 MACHINE LEARNING



Solution: Long Short Term Memory

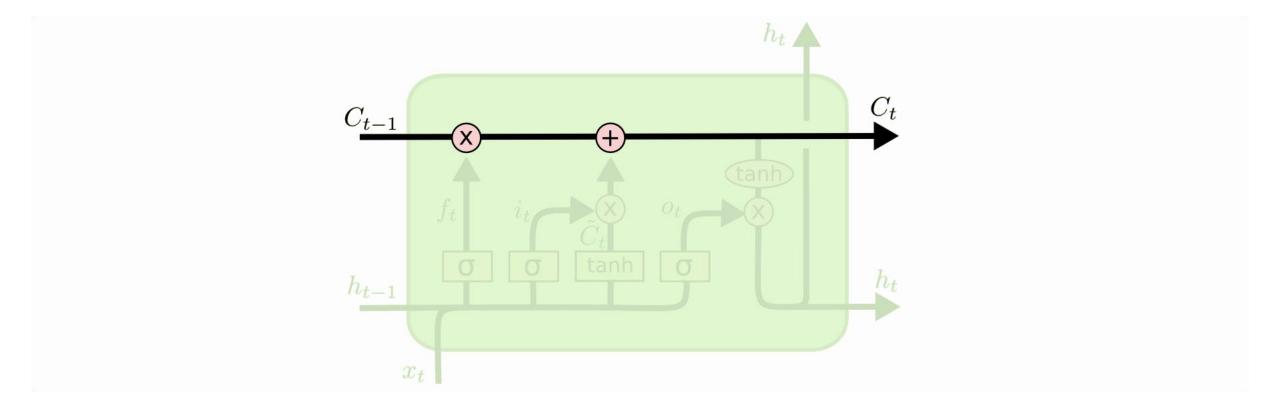
Long Short Term Memory networks (LSTMs) are a special kind of RNN, capable of effectively learning long-term dependencies

- LSTMs are explicitly designed to:
 - Control the vanishing gradient
 - Avoid the long-term dependency problem
 - ▶ Remembering information for long periods of time is practically their default behavior, not something they struggle to learn



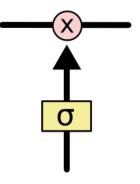


- ▶ The cell state stores an inner information of the chain (inner history)
 - It runs straight down the entire chain, with only minor linear interactions
 - ▶ It's very easy for information to just flow along it unchanged



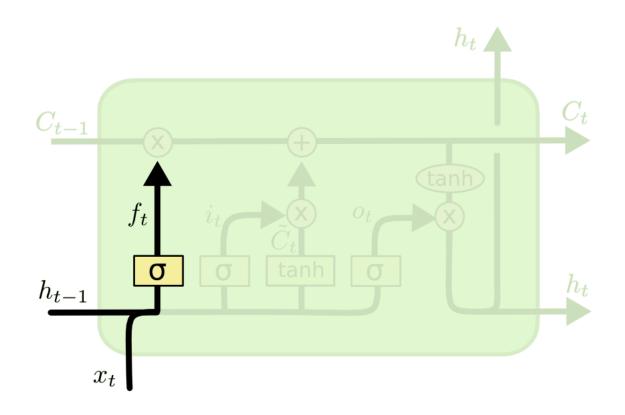


- ▶ The LSTM does have the ability to remove or add information to the state
 - ▶ The information control system is the gate
- ▶ Gates are a way to optionally let information through.
 - ▶ They are composed out of a sigmoid neural net layer and a pointwise multiplication





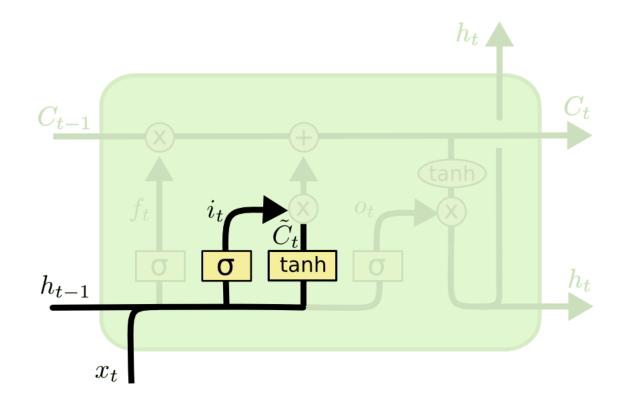
- ▶ The forget gate
 - ▶ It determines how much previous inner history to consider from now on



$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$



- ▶ Input gate and contribution to the current inner history
 - ▶ How much the new input can influence the inner history

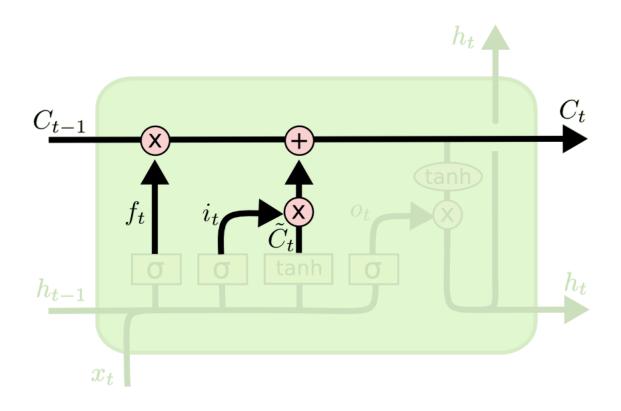


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$



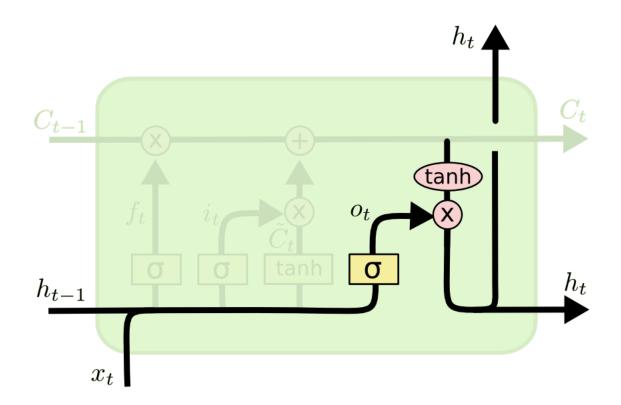
- ▶ Update the inner state
 - ▶ Here LSTMs are tackling the vanishing gradient problem



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



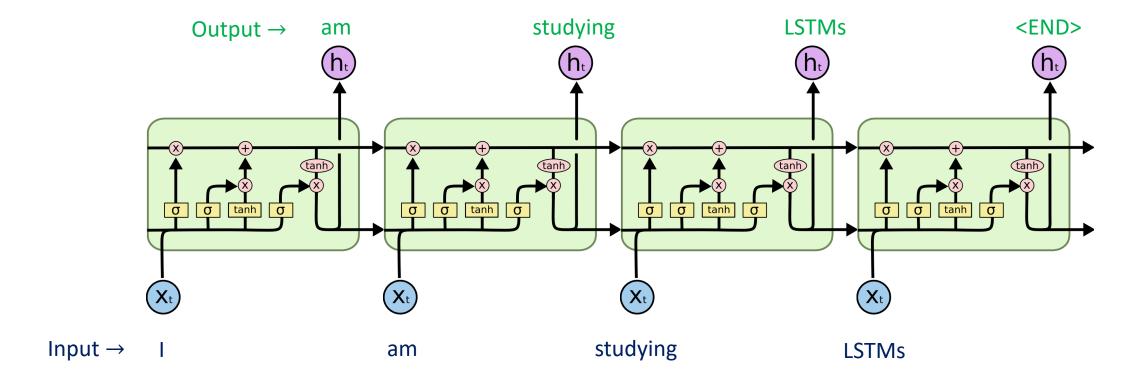
- ▶ Generating the output and output gate
 - ▶ The output is dependent on the cell inner status
 - LSTMs control how much history to use for generating the output



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$
$$h_t = o_t * \tanh(C_t)$$

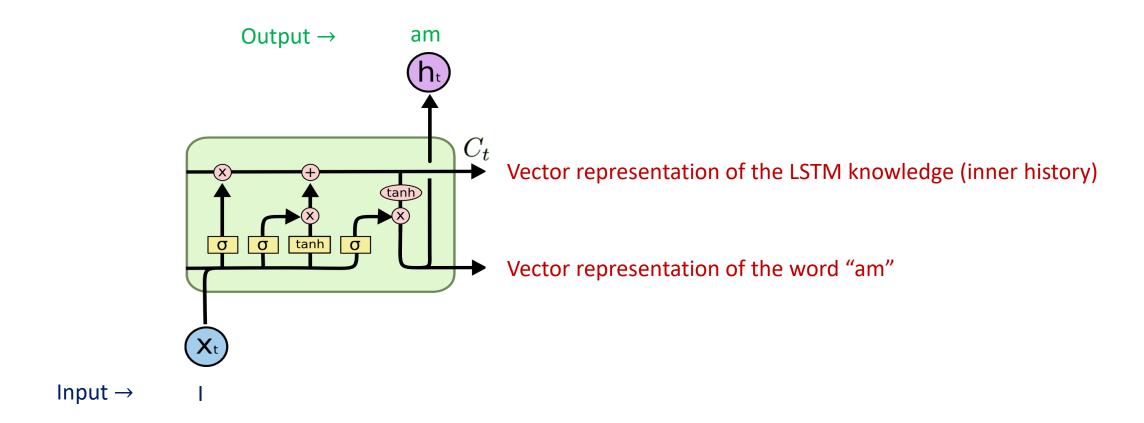
LSTM network: how does it work?

Consider the sentence: "I am studying LSTMs"



2023/2024 **MACHINE LEARNING**

LSTM network: how does it work?



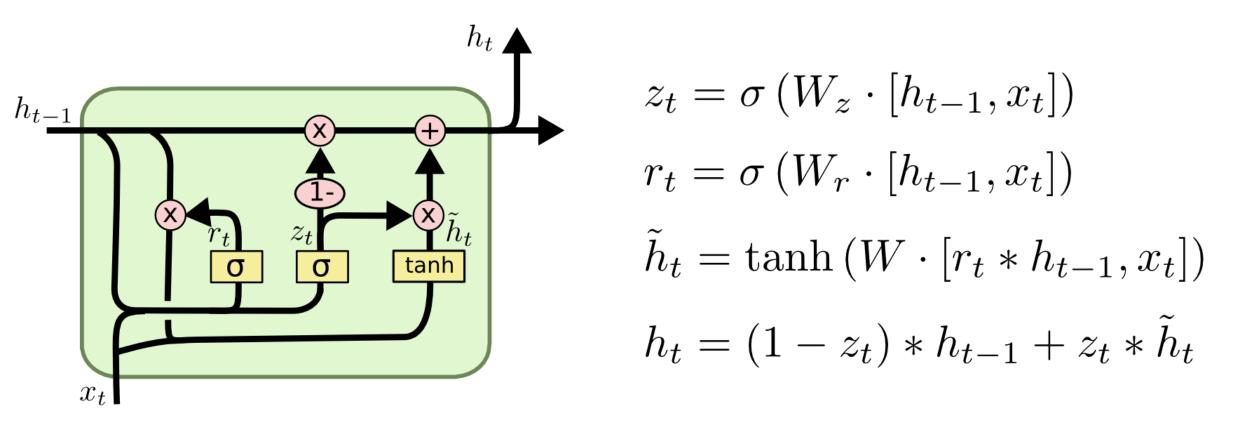
Example

```
from keras.layers import LSTM
model = Sequential()
model.add(Embedding(max features, 32))
model.add(LSTM(32))
model.add(Dense(1, activation='sigmoid'))
model.compile(optimizer='rmsprop',
              loss='binary_crossentropy',
              metrics=['acc'])
history = model.fit(input train, y train,
                    epochs=10,
                    batch size=128,
                    validation split=0.2)
```



Gated Recurrent Unit (GRU)

- Cho, et al. (2014)
- Combine the forget and input gates into a single "update gate"



Setup - Yahoo Finance

```
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

from google.colab import files
files.upload()
apple_training_complete = pd.read_csv(r'AAPL.csv')
apple_training_processed = apple_training_complete.iloc[:, 1:1].values
```

Plot the Dataset

```
plt.figure(figsize=(10,6))
plt.plot(apple training processed, color='blue', label='Apple Stock Price')
plt.xlabel('Date')
plt.ylabel('Apple Stock Price')
plt.legend()
plt.show()
                                              Apple Stock Price
                                        350
                                        325
                                      Apple Stock Price
                                        300
                                       275
                                        250
                                        225
                                        200
                                                                 100
                                                                           150
                                                                                     200
                                                                                                250
                                                       50
                                                                      Date
```

Preprocessing

Building the Network

```
from keras.models import Sequential
from keras.layers import Dense
from keras.layers import LSTM
from keras.layers import Dropout
model = Sequential()
model.add(LSTM(units=50, return sequences=True, input shape=(features set.shape[1], 1)))
model.add(Dropout(0.2))
model.add(LSTM(units=50, return sequences=True))
model.add(Dropout(0.2))
model.add(LSTM(units=50, return sequences=True))
model.add(Dropout(0.2))
model.add(LSTM(units=50))
model.add(Dropout(0.2))
model.add(Dense(units = 1))
model.compile(optimizer = 'adam', loss = 'mean squared error')
```



```
model.fit(features_set, labels, epochs = 100, batch_size = 32)
```



```
from google.colab import files
files.upload()
apple testing complete = pd.read csv(r'AAPL-Test.csv')
apple testing processed = apple testing complete.iloc[:, 1:1].values
apple_total = pd.concat((apple_training_complete['Open'], apple testing complete['Open']), axis=0)
test inputs = apple total[len(apple total) - len(apple testing complete) - 60:].values
test inputs = test inputs.reshape(-1,1)
test inputs = scaler.transform(test inputs)
test features = []
for i in range(60, 80):
   test features.append(test inputs[i-60:i, 0])
test features = np.array(test features)
test features = np.reshape(test features, (test features.shape[0], test features.shape[1], 1))
```



