

LSTM AND GRU

Vincent Barra

LIMOS, UMR 6158 CNRS, Université Clermont Auvergne

Sequence modeling

- ▶ Handle variable-length inputs
- ▶ Share parameters across the sequence
- ▶

Sequence modeling

- ▶ Handle variable-length inputs
- ▶ Share parameters across the sequence
- ▶ **Keep track of long-term dependencies**

Paul is a small boy and is hungry. **He** goes to the restaurant and eats a lot.

Sequence modeling

- ▶ Handle variable-length inputs
- ▶ Share parameters across the sequence
- ▶ **Keep track of long-term dependencies**

Paul is a small boy and is hungry. **He** goes to the restaurant and eats a lot.

Solutions

- ▶ Activation function: using ReLU prevents σ' from shrinking the gradients when $x > 0$
- ▶ Weights and bias initialisation can help ($\mathbf{W} = \mathbf{I}$, $\mathbf{b}=0$)
- ▶ **Think of another architecture**

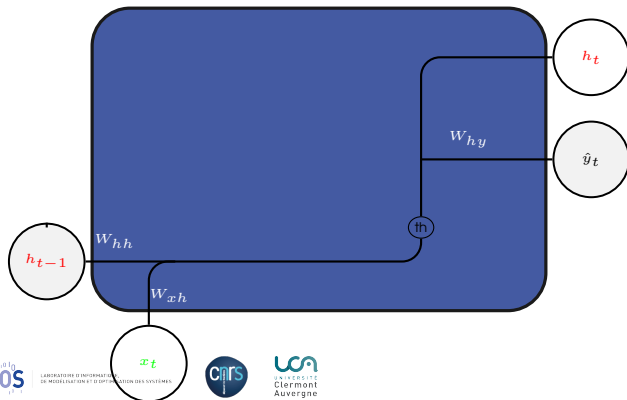
RECALL: RNN

- Update hidden state

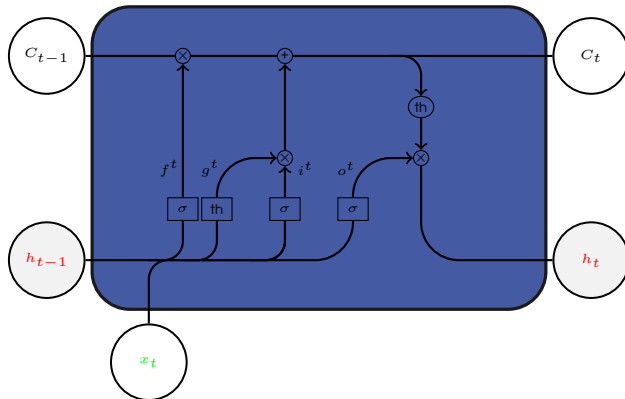
$$h_t = \tanh(W_{xh}^T x_t + W_{hh}^T h_{t-1})$$

- Compute output vector

$$\hat{y}_t = W_{hy}^T h_t$$



LSTM



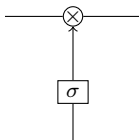
LSTM

Properties

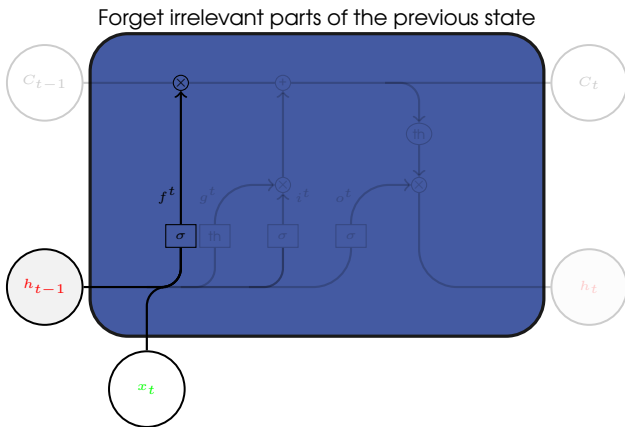
- ▶ Long Short Term Memory networks rely on gated cell to allow information tracking through time, even for several timesteps
- ▶ Contains computational blocks controlling information flow
- ▶ Key concept: a persistent cell state C_t module, representation of past history
- ▶ in Keras: `tensorflow.keras.layers.LSTM(num_units)`

Gates

- ▶ Information is added or removed using gate structures.
- ▶ can let information through using eg. sigmoid and pointwise multiplication



INSIDE LSTM: FORGET GATE



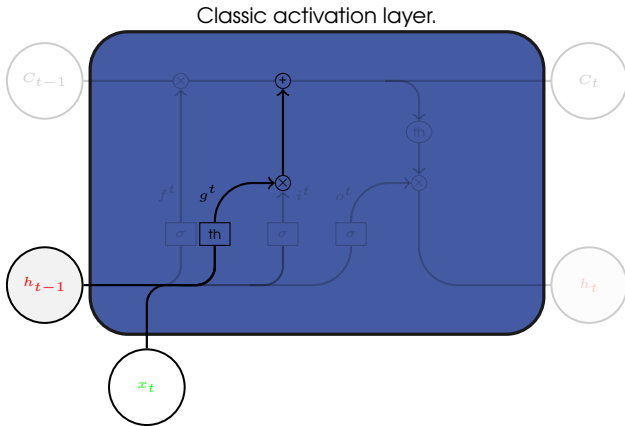
INSIDE LSTM: FORGET GATE

$$f^t = \sigma \left(\mathbf{W}_f^\top [\mathbf{x}_t, \mathbf{h}_{t-1}] + b_f \right)$$

example: the state keeps the gender of the hero in a text, for a good use of pronouns. If a new hero appears, forget the gender.

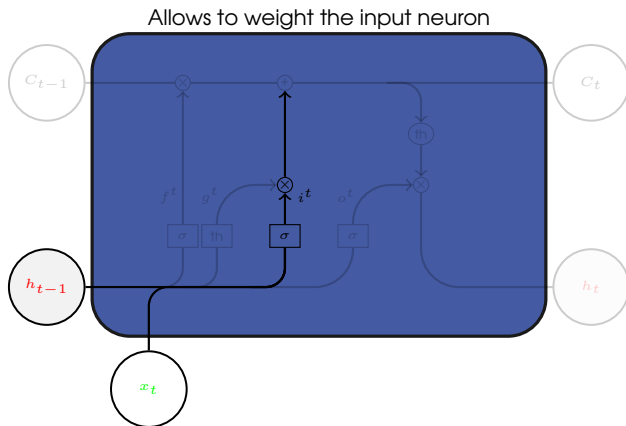
- ▶ $f_t = 0$: completely forgetting previous state
- ▶ $f_t = 1$: completely keeping previous state
- ▶ b_f should be initialized with large values so that initially $f_t \approx 1$

INSIDE LSTM: INPUT NEURON



$$g^t = \tanh(\mathbf{W}_C^T [\mathbf{x}_t, \mathbf{h}_{t-1}] + b_C)$$

INSIDE LSTM: UPDATE GATE

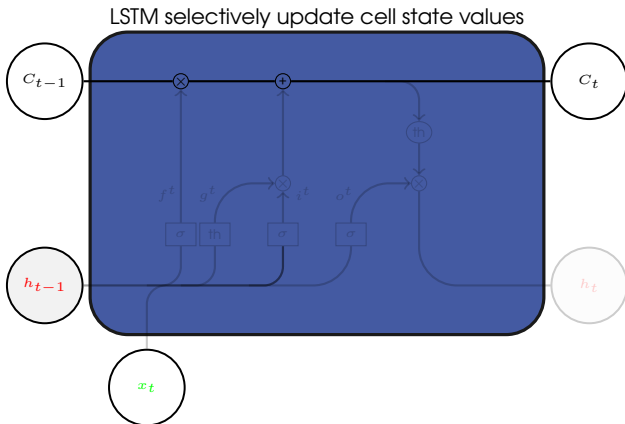


INSIDE LSTM: UPDATE GATE

$$i^t = \sigma \left(\mathbf{W}_i^\top [\mathbf{x}_t, h_{t-1}] + b_i \right)$$

- ▶ Decide the importance of the contribution of the input neuron at each timestep.
- ▶ g^t and i^t store relevant new information in the current state.

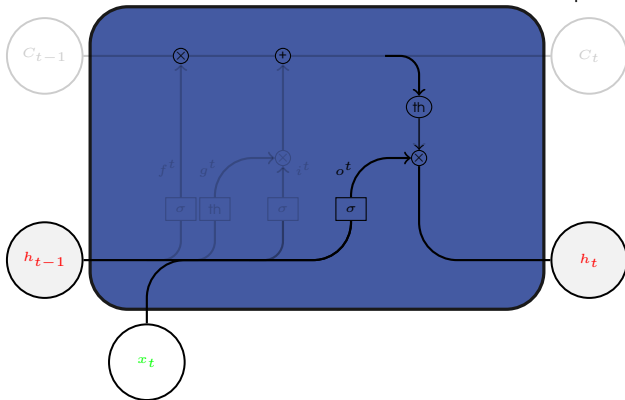
INSIDE LSTM: UPDATE STEP



$C_t = g^t \cdot i^t + C_{t-1} \cdot f^t$: assures that the derivatives of the loss w.r.t. C_t does not vanish.

INSIDE LSTM: OUTPUT GATE

Controls what information is sent to the next timestep



$$o^t = \sigma \left(W_o^T [\mathbf{x}_t, \mathbf{h}_{t-1}] + b_o \right) \quad h_t = o^t \tanh(C_t)$$

In summary:

$$f^t = \sigma \left(\mathbf{W}_f^\top [\mathbf{x}_t, \mathbf{h}_{t-1}] + b_f \right)$$

$$g^t = \tanh \left(\mathbf{W}_C^\top [\mathbf{x}_t, \mathbf{h}_{t-1}] + b_C \right)$$

$$i^t = \sigma \left(\mathbf{W}_i^\top [\mathbf{x}_t, \mathbf{h}_{t-1}] + b_i \right)$$

$$C_t = g^t \cdot i^t + C_{t-1} \cdot f^t$$

$$o^t = \sigma \left(\mathbf{W}_o^\top [\mathbf{x}_t, \mathbf{h}_{t-1}] + b_o \right)$$

$$\mathbf{h}_t = o^t \tanh(C_t)$$

KEY CONCEPTS

Concepts

- 1 Maintain a separate state C_t from what is outputted
- 2 Use gates to control the information flow
 - ▶ can forget information (f^t)
 - ▶ can store relevant information from the x_t (g^t)
 - ▶ can selectively update state (i^t)
 - ▶ can return a filtered version of the state (o^t)
- 3 introduction of self-loops to produce paths where gradients can flow for long durations
- 4 Since the backward flow from C_t to C_{t-1} is direct, backpropagation through time is computed with uninterrupted gradient flow.

EXAMPLE

- ▶ Task: predicting the next word of a sentence based on previous ones.
- ▶ Hypothesis: the cell state C_t contains the gender of the subject

Paul is a clever guy. He has his own car. Cathy and Mary are their sisters.

- 1 First step (f^t): C_{t-1} contains the gender of the subject (Paul) to use proper pronouns (his). If a new subject arrive (Cathy and Mary), we may want to forget the old gender.
- 2 Second step: (g^t, i^t): what kind of information we want to store in C_t ? Here we may want to add the gender of the new subject (Cathy and Mary) to C_t
- 3 Third step (o^t): what we want to output. Since we see a new subject, we may want to output information relevant to verb (are)(eg. singular or plural subject)

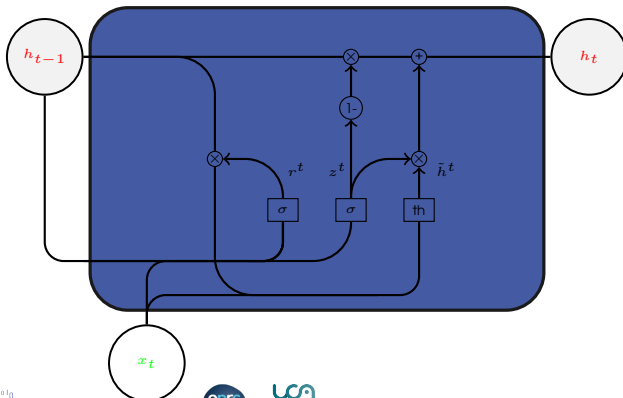
VARIANTS

- ▶ LSTM with peephole connections
- ▶ Gates have access to C_{t-1}
- ▶ Bi-directional recurrent networks
- ▶ Gated Recurrent Units (GRU)
- ▶ Skip LSTMs
- ▶ ...

GATED RECURRENT UNITS

Simplified version of LSTM:

- ▶ no o^t neither C_t
- ▶ only two gates
- ▶ easier to train



GATED RECURRENT UNITS

- ▶ reset gate r^t : determines how to combine \mathbf{x}_t with h_{t-1}
- ▶ update gate z^t : what quantity of memory must be preserved (like a combination of f^t and g^t)

$$r^t = \sigma(\mathbf{W}_r^\top [\mathbf{x}_t, h_{t-1}] + b_r)$$

$$z^t = \sigma(\mathbf{W}_z^\top [\mathbf{x}_t, h_{t-1}] + b_z)$$

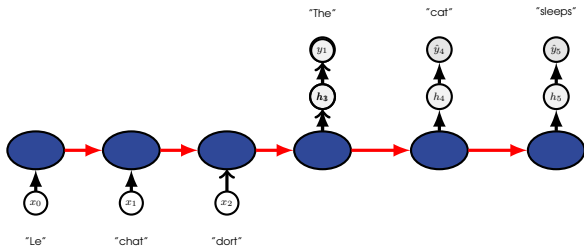
$$\tilde{h}^t = \tanh(\mathbf{W}^\top [\mathbf{x}_t, r^t h_{t-1}] + b_h)$$

$$h_t = (1 - z^t) h_{t-1} + z^t \tilde{h}^t$$

Gating network signals control how the present input and previous memory are used to update the current activation and produce the current state.

MACHINE TRANSLATION

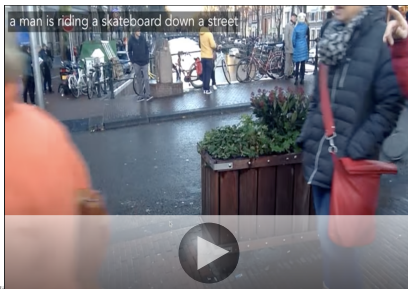
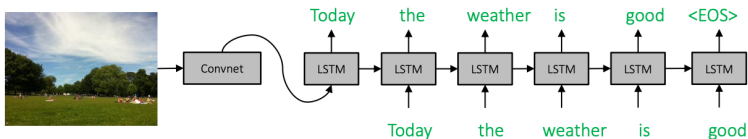
Many-to-Many architecture



- ▶ the sentence in the source language is a sequence
- ▶ it is transformed in a latent space by an LSTM encoder
- ▶ the sentence in the target language is a sequence
- ▶ it is captured by a LSTM decoder

IMAGE CAPTIONING

Many-to-Many architecture



TIMESERIES MODELING

Environmental modeling

KERAS

```

model = Sequential()
# 1 layer LSTM with batch size = 128
model.add(LSTM(128, input_shape=(...)))
model.add(Dense(..))
# Predicting the next word
model.add(Activation('softmax'))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
print(model.summary())

# Using embeddings
embedding_size = 32
model = Sequential()
model.add(Embedding(top_words, embedding_size, input_length=max_length))
model.add(LSTM(100))
model.add(Dense(1, activation='sigmoid'))
model.compile(loss='binary_crossentropy', optimizer='adam', metrics=['accuracy'])
print(model.summary())

```