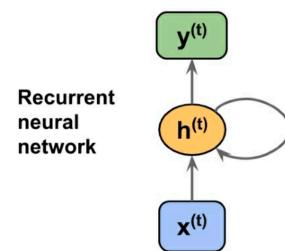
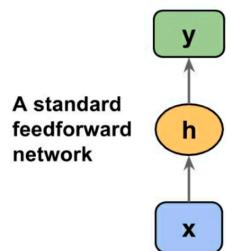


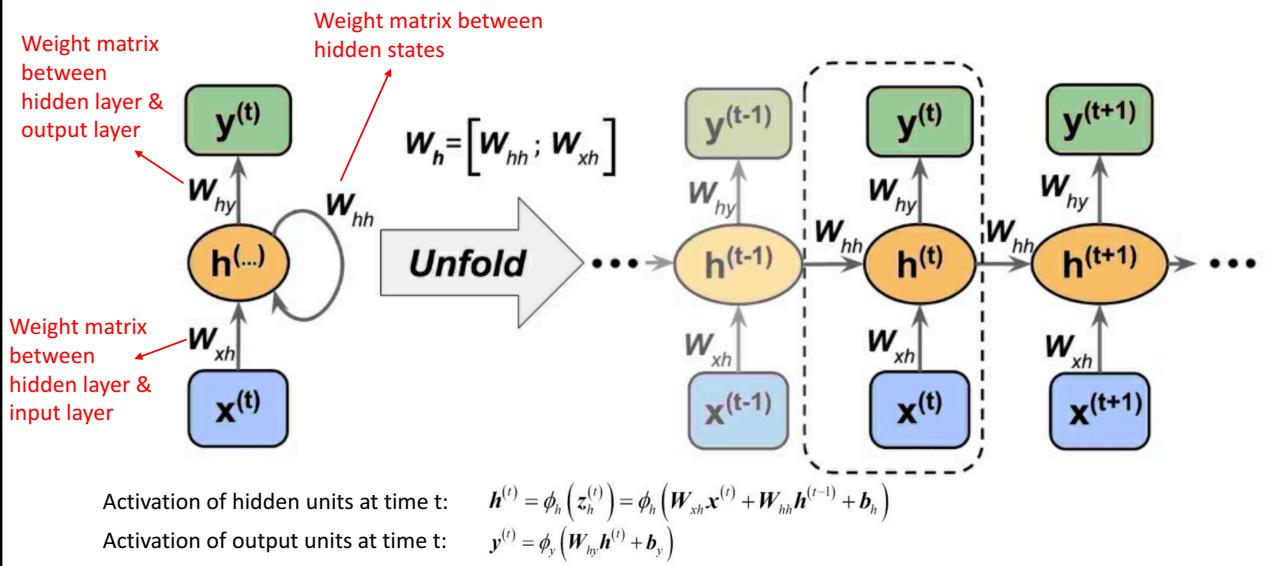
Introduction of LSTM

Neural network & RNN

- IID data
- Neural network
- Time dependent data
- RNN



Unfold of RNN



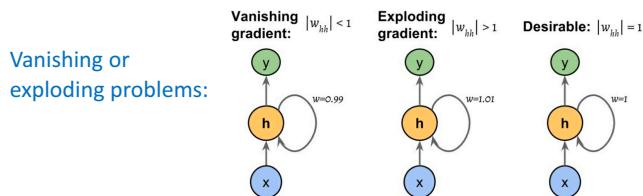
Training RNNs

Overall loss function: $L = \sum_{t=1}^T L^{(t)}$ Sum of all loss functions at times t=1 to t=T

Gradient of L w.r.t \mathbf{W}_{hh} : $\frac{\partial L^{(t)}}{\partial \mathbf{W}_{hh}} = \frac{\partial L^{(t)}}{\partial \mathbf{y}^{(t)}} \times \frac{\partial \mathbf{y}^{(t)}}{\partial \mathbf{h}^{(t)}} \times \left(\sum_{k=1}^t \frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} \times \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{W}_{hh}} \right)$

Recall: activation of hidden states $\mathbf{h}^{(t)} = \phi_h(\mathbf{z}_h^{(t)}) = \phi_h(W_{xh}\mathbf{x}^{(t)} + W_{hh}\mathbf{h}^{(t-1)} + \mathbf{b}_h)$

Recurrent: multiplication of adjacent steps $\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} = \prod_{i=k+1}^t \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{h}^{(i-1)}}$



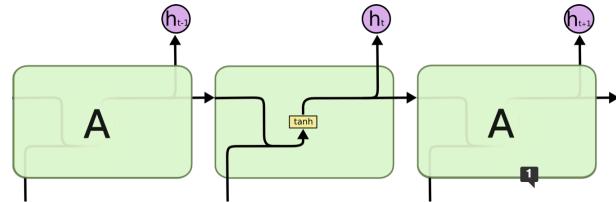
Solutions:

- Truncated backpropagation through time (TBPTT)
- Long short-term memory (LSTM)

Structure of LSTM

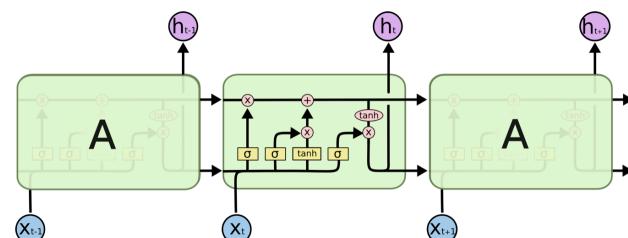
RNN

$$\begin{aligned} h^{(t)} &= \phi_h(z_h^{(t)}) = \phi_h(W_{xh}x^{(t)} + W_{hh}h^{(t-1)} + b_h) \\ y^{(t)} &= \phi_y(W_{hy}h^{(t)} + b_y) \end{aligned}$$



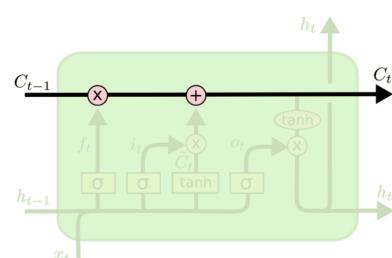
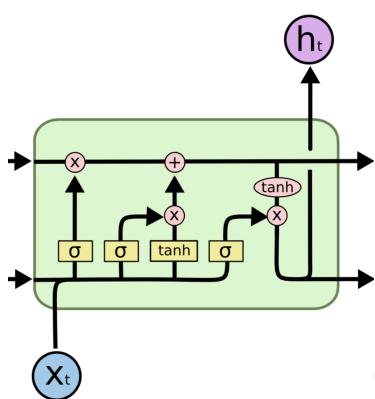
LSTM

- ❖ Building block: memory cell
- ❖ Recurrent edge with w=1
- ❖ Value associated with recurrent edge: cell state

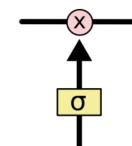
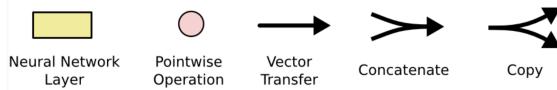


Elements of LSTM

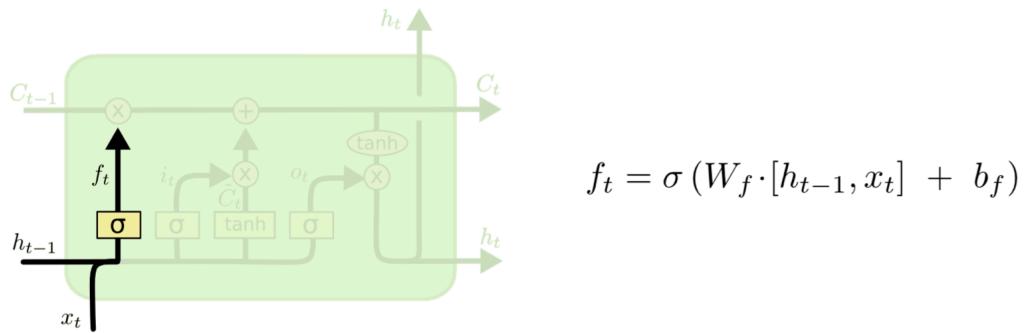
Core: cell state
Can add or remove information regulated by gates



Gate: a way to optionally let information through

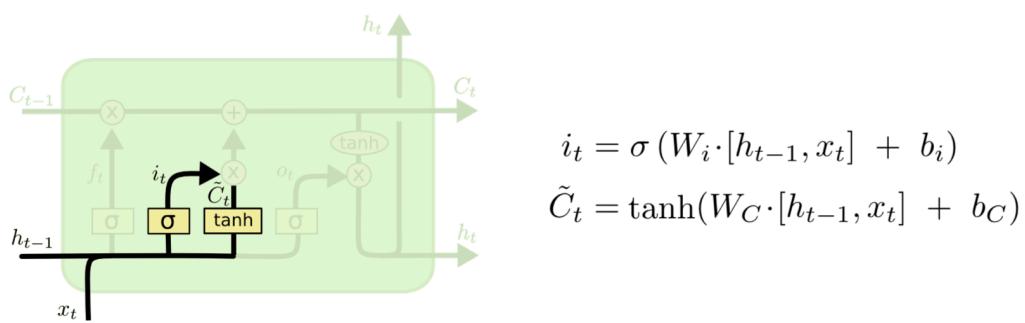


Step-by-step LSTM Walk Through

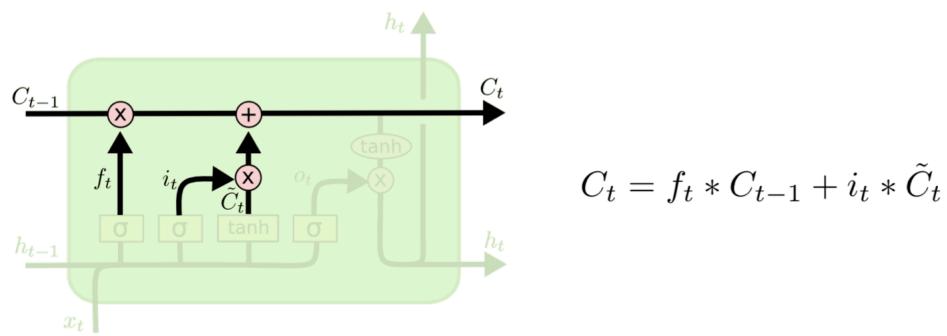


Decide what information to throw away from the cell state
Forget gate layer: a sigmoid

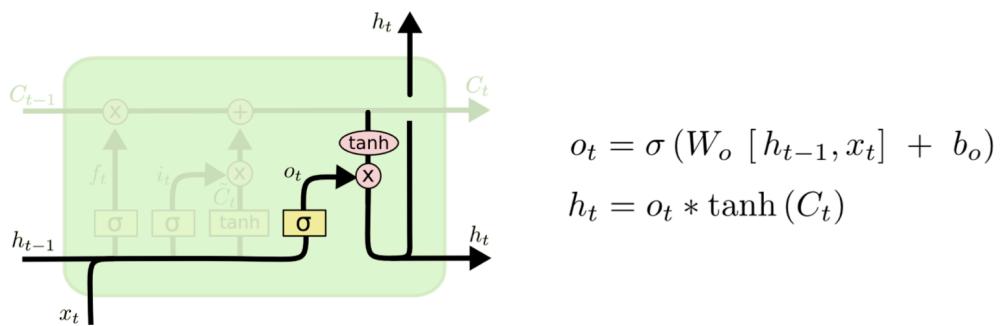
Step-by-step LSTM Walk Through



Step-by-step LSTM Walk Through



Step-by-step LSTM Walk Through



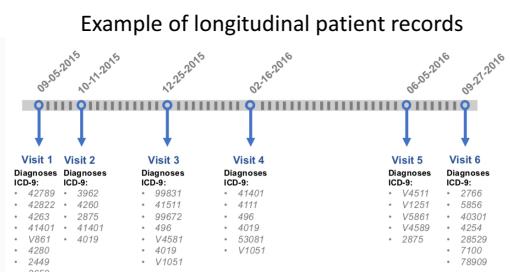
Challenges of using LSTM in medical data

- Time intervals of each visit are not evenly distributed
 - T-LSTM
- Poor interpretability
 - Attention mechanism

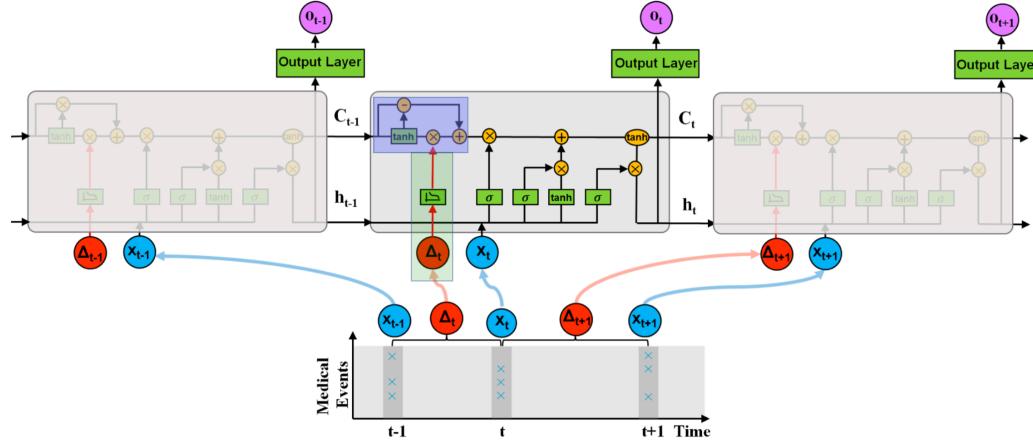
Time-aware LSTM

- Time irregularity is important in clinical data

- Variation in time gap could be indicative of certain impending disease condition
- Dependency on the previous memory should not play an active role in prediction of current status if the time interval is too large

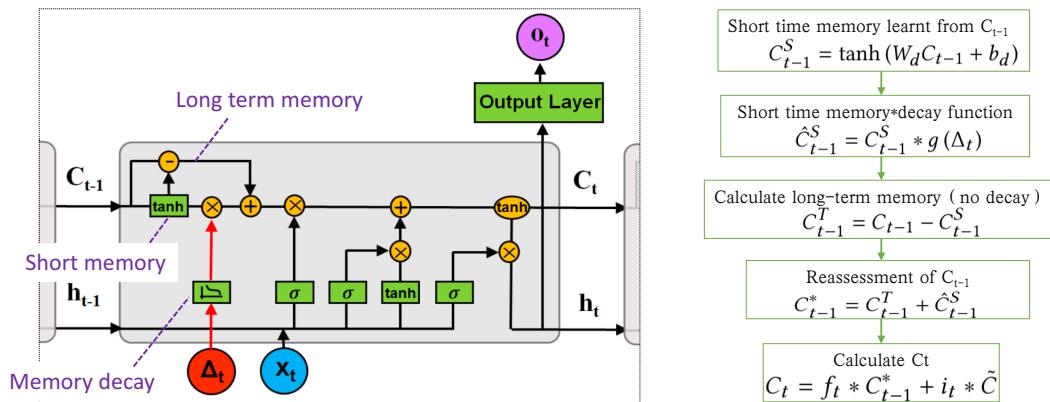


Unit of T-LSTM



Difference compared with standard LSTM:
memory decomposition + decay of short memory

Algorithm of T-LSTM



Trade accuracy for interpretability in traditional models

- Identifying a set of rules
 - e.g. via decision trees
- Case-based reasoning by finding similar patients
 - e.g. via k-nearest neighbors & distance metric learning
- Identifying a list of risk factors
 - e.g. via LASSO coefficients

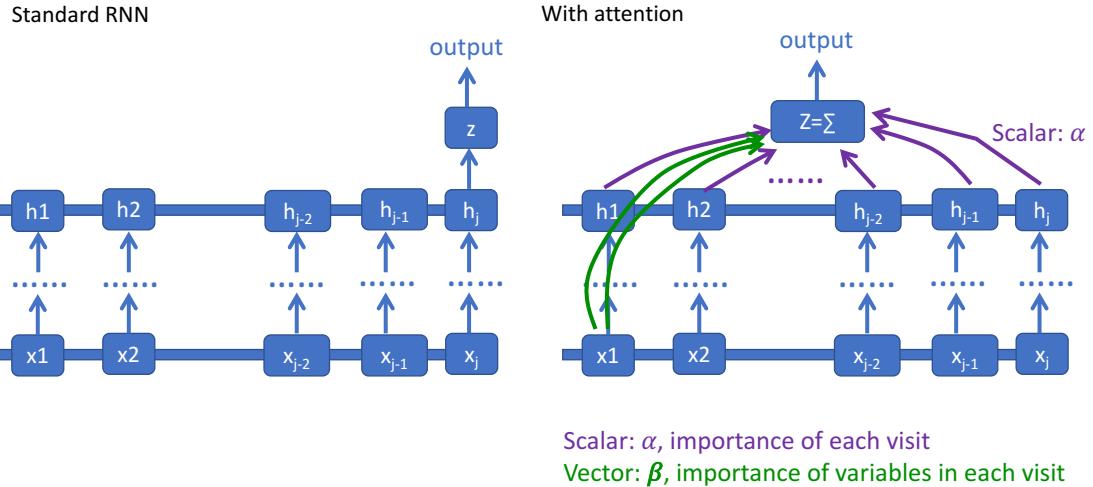
Attention mechanism

Figure 5. Examples of mistakes where we can use attention to gain intuition into what the model saw.

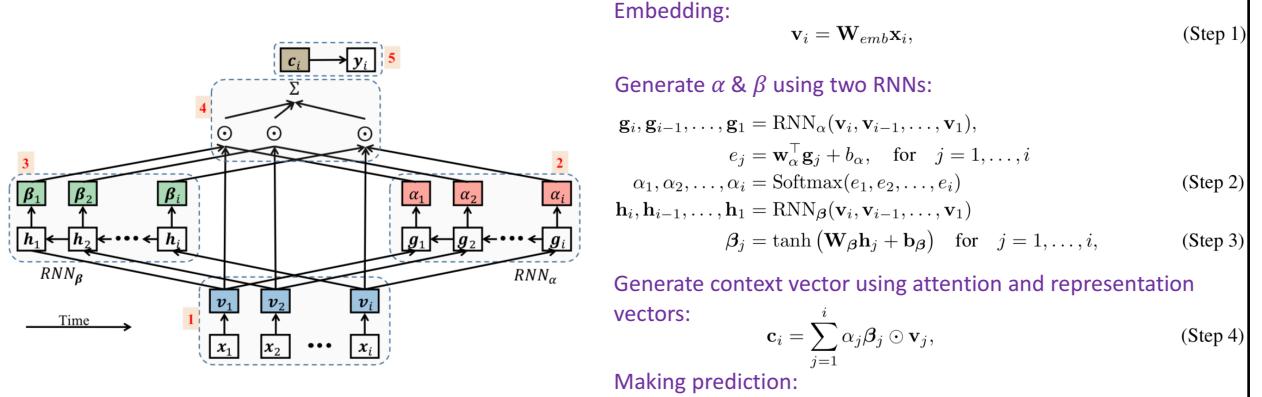


- Mimic working mechanism of human brain
- Used in figure caption, image classification, language translation
- Helps to assign different weights to input

Introducing attention mechanism in LSTM



RETAIN: reverse time attention model



Interpreting RETAIN:

$$p(\mathbf{y}_i | \mathbf{x}_1, \dots, \mathbf{x}_i) = p(\mathbf{y}_i | \mathbf{c}_i) = \text{Softmax}(\mathbf{W}\mathbf{c}_i + \mathbf{b})$$

$$p(\mathbf{y}_i | \mathbf{x}_1, \dots, \mathbf{x}_i) = p(\mathbf{y}_i | \mathbf{c}_i) = \text{Softmax}\left(\mathbf{W}\left(\sum_{j=1}^i \alpha_j \boldsymbol{\beta}_j \odot \mathbf{v}_j\right) + \mathbf{b}\right)$$

$$\begin{aligned} p(\mathbf{y}_i | \mathbf{x}_1, \dots, \mathbf{x}_i) &= \text{Softmax}\left(\mathbf{W}\left(\sum_{j=1}^i \alpha_j \boldsymbol{\beta}_j \odot \sum_{k=1}^r x_{j,k} \mathbf{W}_{emb}[:, k]\right) + \mathbf{b}\right) \\ &= \text{Softmax}\left(\sum_{j=1}^i \sum_{k=1}^r x_{j,k} \alpha_j \mathbf{W}(\boldsymbol{\beta}_j \odot \mathbf{W}_{emb}[:, k]) + \mathbf{b}\right) \end{aligned}$$

$$\omega(\mathbf{y}_i, x_{j,k}) = \underbrace{\alpha_j \mathbf{W}(\boldsymbol{\beta}_j \odot \mathbf{W}_{emb}[:, k])}_{\text{Contribution coefficient}} \underbrace{x_{j,k}}_{\text{Input value}},$$

references

- Background LSTM: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>
- T-LSTM:
 - Original paper:
http://biometrics.cse.msu.edu/Publications/MachineLearning/Baytasetal_PatientSubtypingViaTimeAwareLSTMNetworks.pdf
 - Github:
<https://github.com/illidanlab/T-LSTM>
- RETAIN:
 - Original paper:
<https://arxiv.org/abs/1608.05745>
 - Github:
<https://github.com/mp2893/retain>