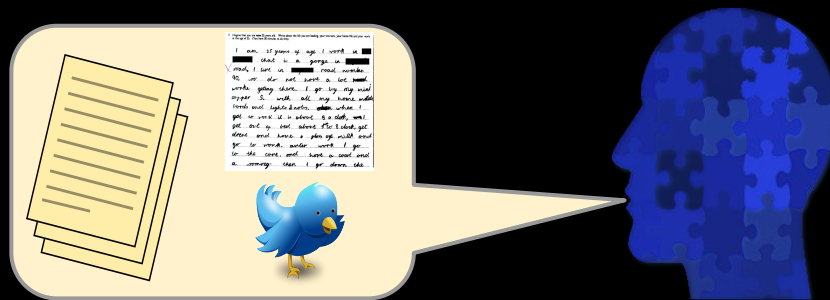


Recurrent Neural Networks for Language Modeling

CSE354 - Spring 2021
Natural Language Processing

Tasks



- Language Modeling:
Generate next word, sentence

how?



≈ capture hidden

representation of sentences.

- Word, Document Classification
(named entity tagging; sentiment
analysis using sequence, etc...)

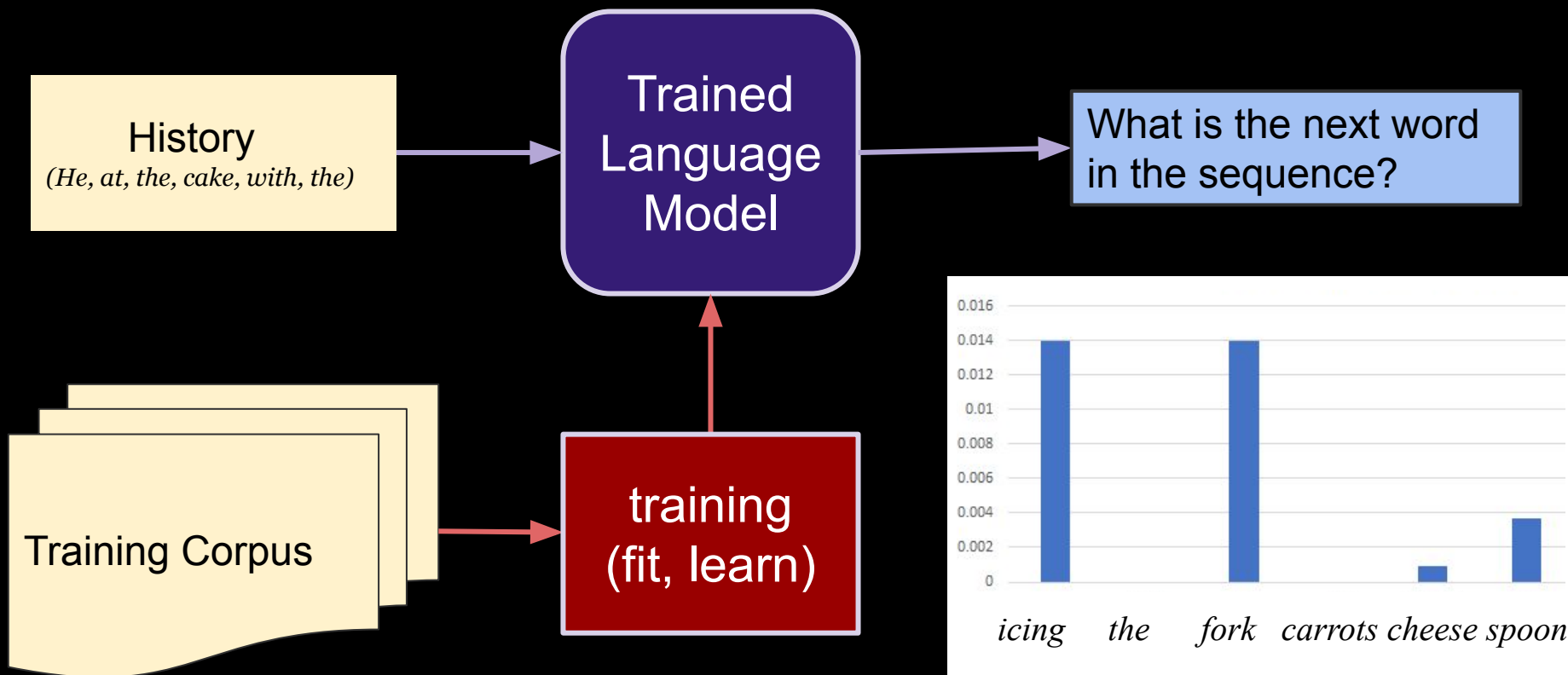
- Recurrent Neural Network and
Sequence Models

Language Modeling

Task: Estimate $P(w_n | w_1, w_2, \dots, w_{n-1})$
:probability of a next word given history
 $P(\text{fork} | \text{He ate the cake with the}) = ?$

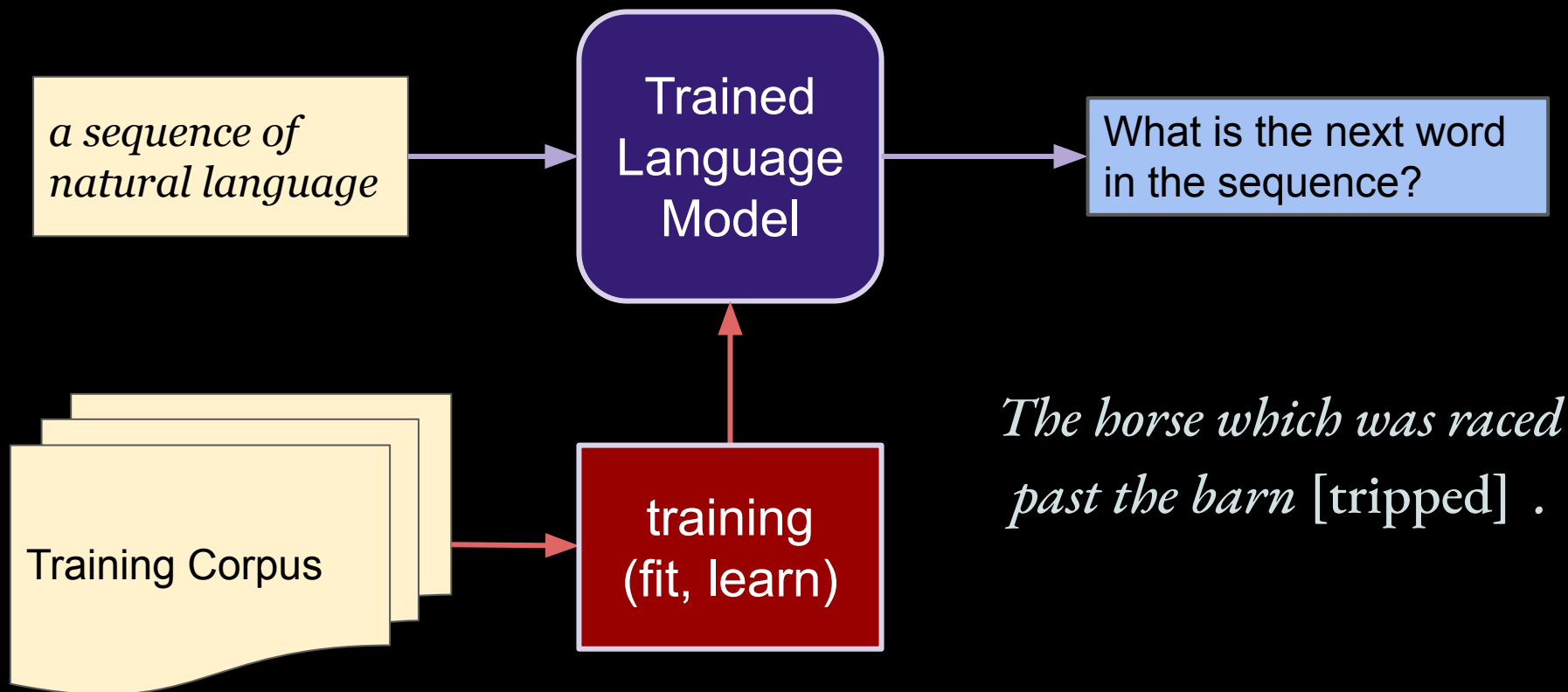
Language Modeling

Task: Estimate $P(w_n | w_1, w_2, \dots, w_{n-1})$
:probability of a next word given history
 $P(\text{fork} | \text{He ate the cake with the}) = ?$



Language Modeling

Building a model (or system / API) that can answer the following:

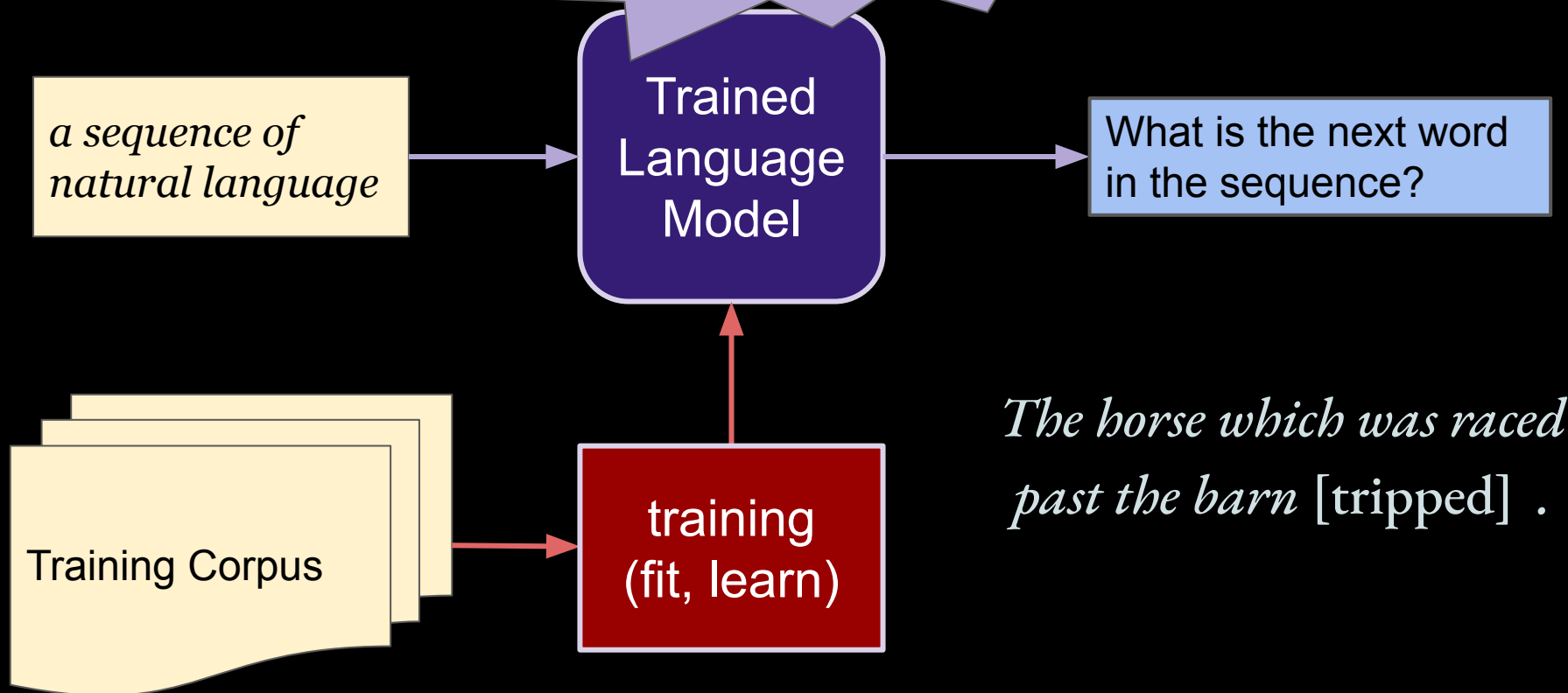


Language Modeling

Building a model (or s

To fully capture natural language, models get very complex!

er the following:



*The horse which was raced
past the barn [tripped] .*

Neural Networks: Graphs of Operations (excluding the optimization nodes)

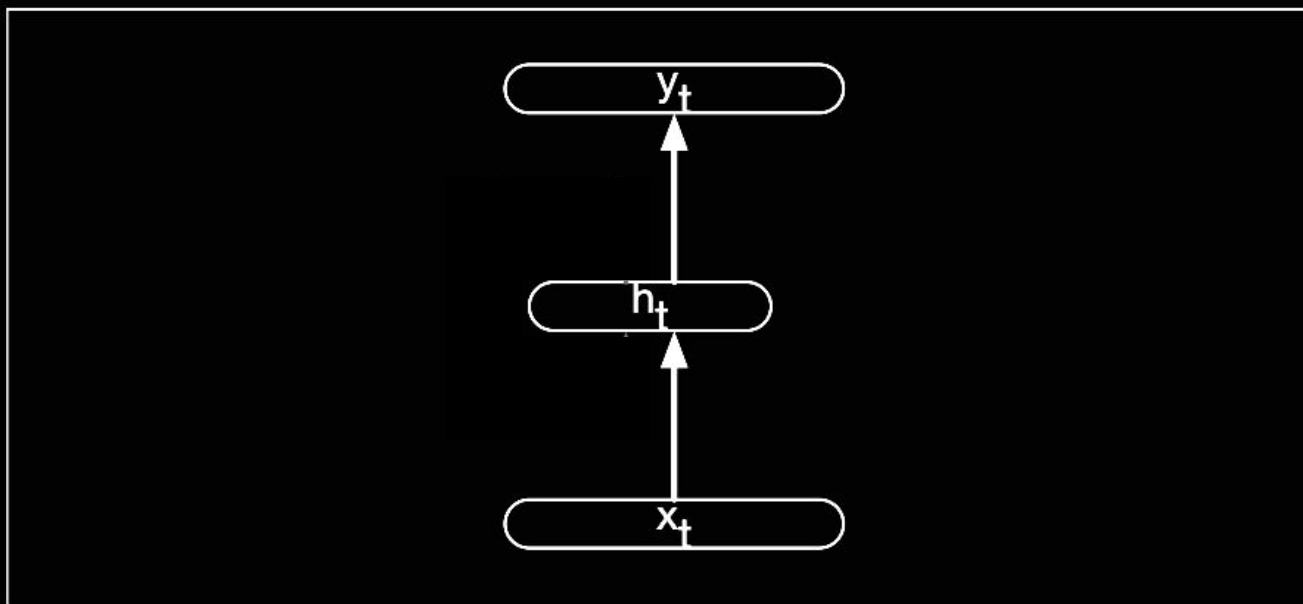


Figure 9.2 Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

Neural Networks: Graphs of Operations (excluding the optimization nodes)

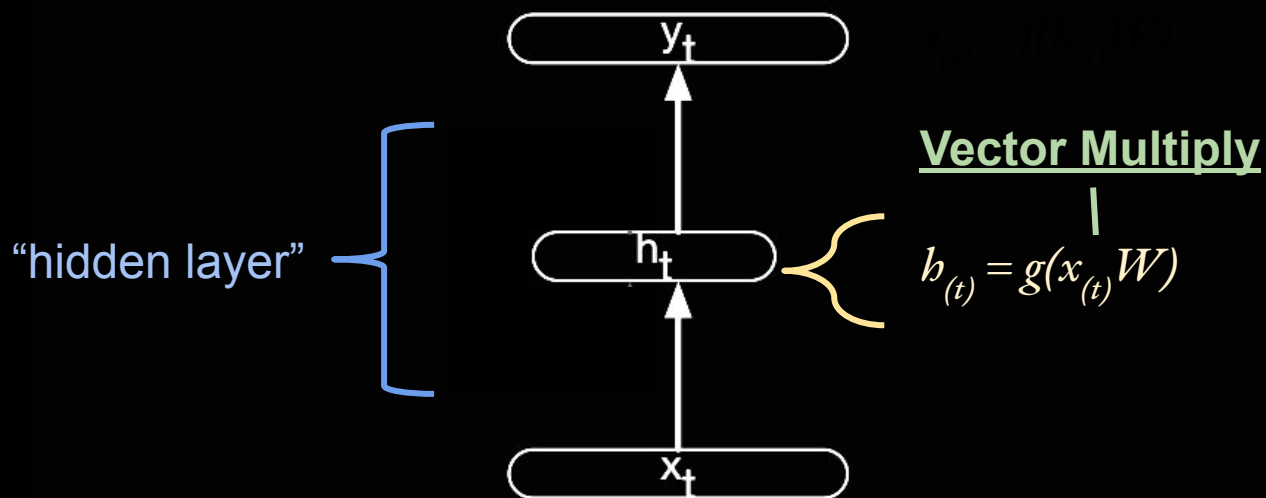


Figure 9.2 Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

Neural Networks: Graphs of Operations (excluding the optimization nodes)

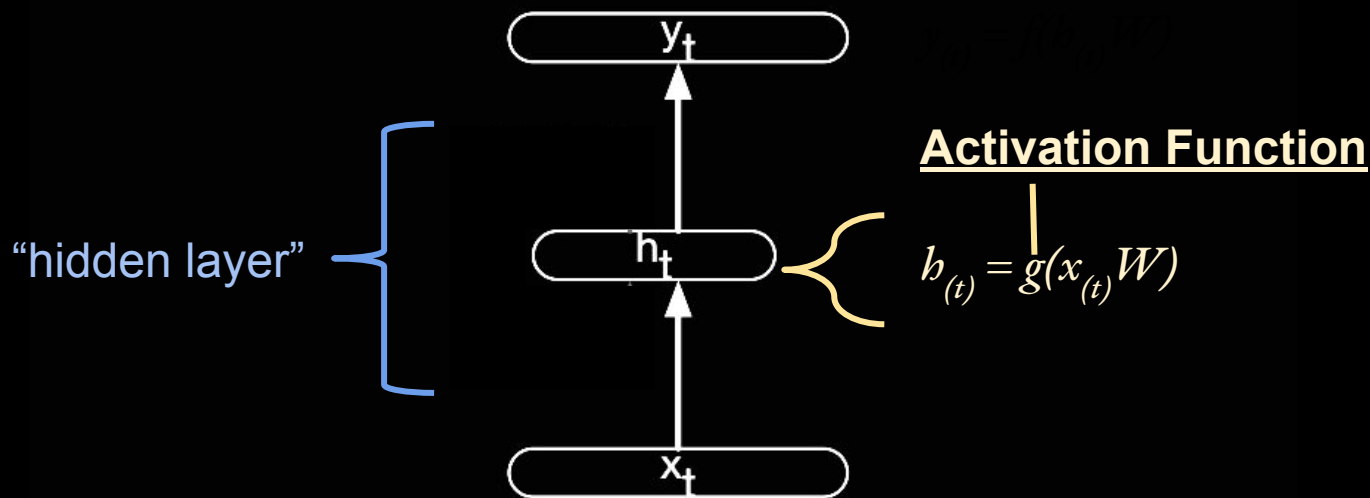
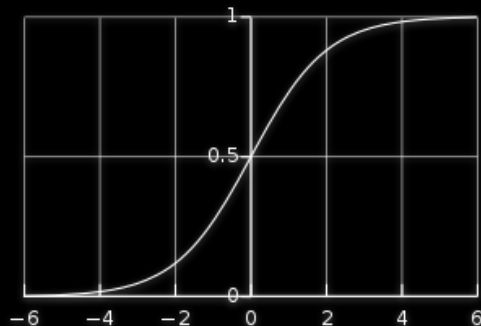


Figure 9.2 Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

Common Activation Functions

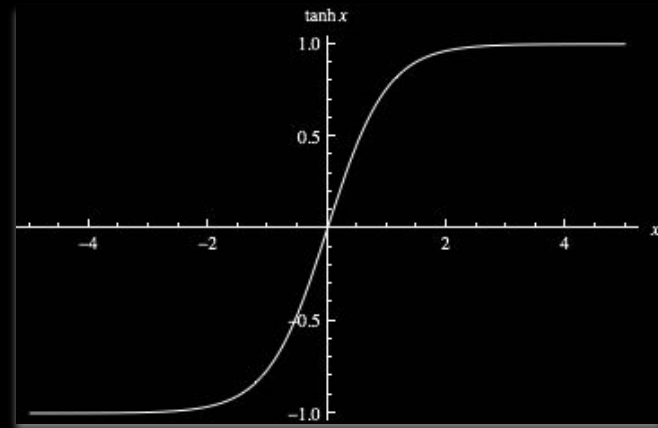
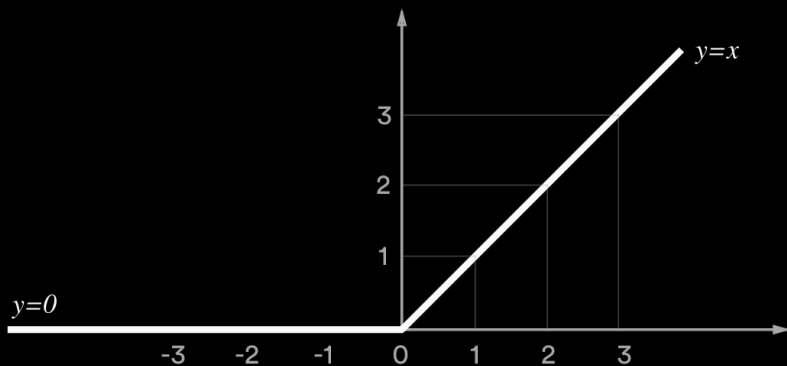
$$z = b_{(t)}W$$

Logistic: $\sigma(z) = 1 / (1 + e^{-z})$

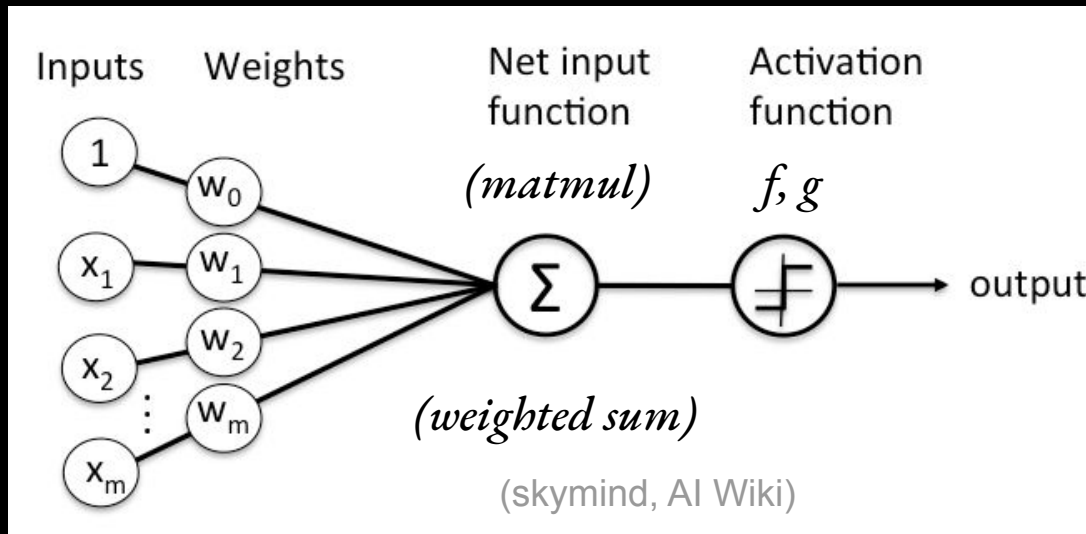


Hyperbolic tangent: $\tanh(z) = 2\sigma(2z) - 1 = (e^{2z} - 1) / (e^{2z} + 1)$

Rectified linear unit (ReLU): $ReLU(z) = \max(0, z)$



Neural Networks: Graphs of Operations (excluding the optimization nodes)



Activation Function

$$h_{(t)} = g(x_{(t)}W)$$

Figure 9.2 Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep.

(Jurafsky, 2019)

Neural Networks: Graphs of Operations

(excluding the optimization nodes)

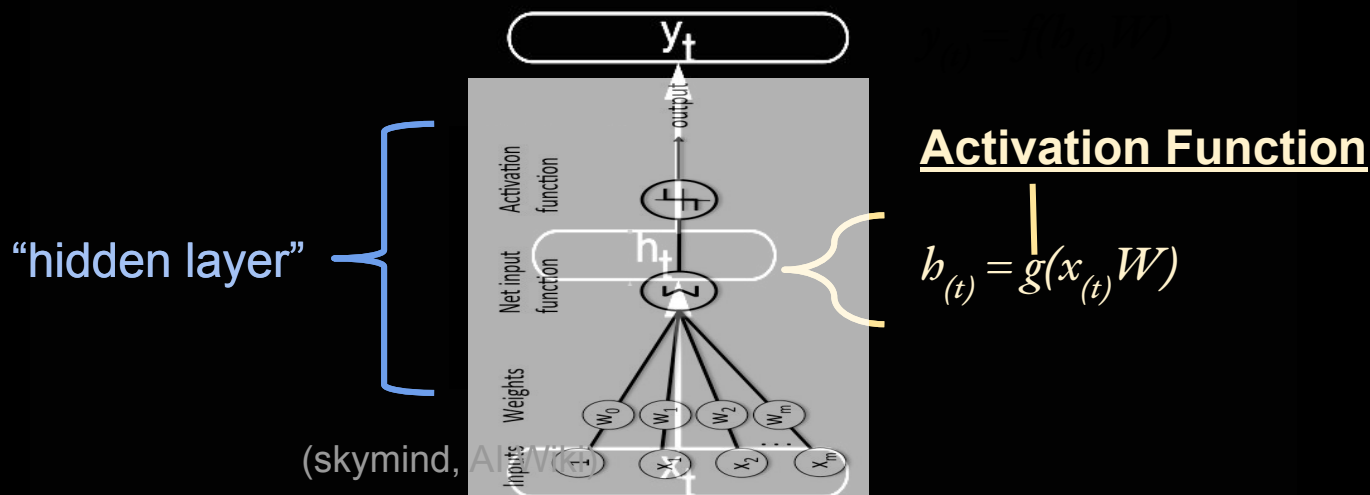


Figure 9.2 Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

Neural Networks: Graphs of Operations (excluding the optimization nodes)

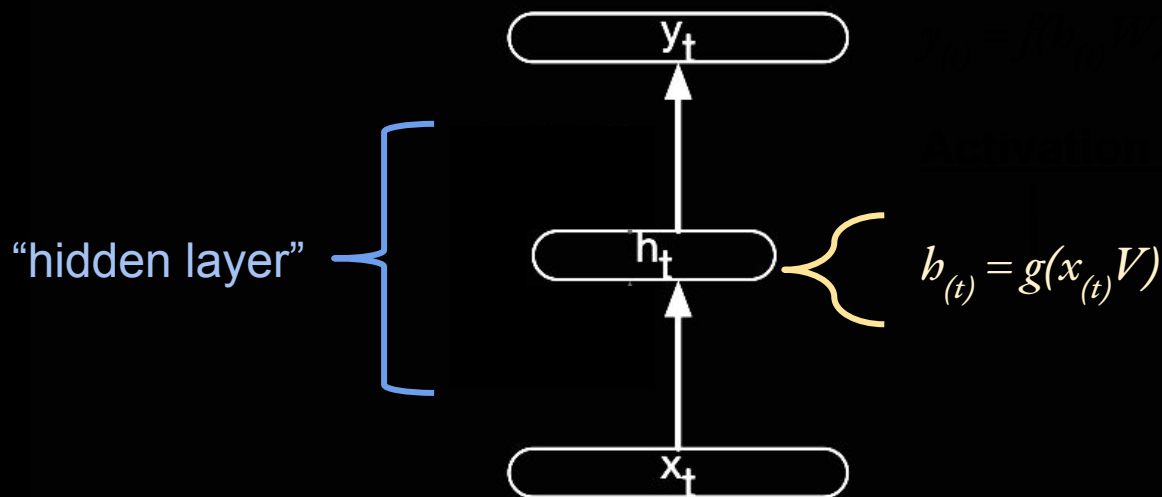


Figure 9.2 Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

Neural Networks: Graphs of Operations (excluding the optimization nodes)

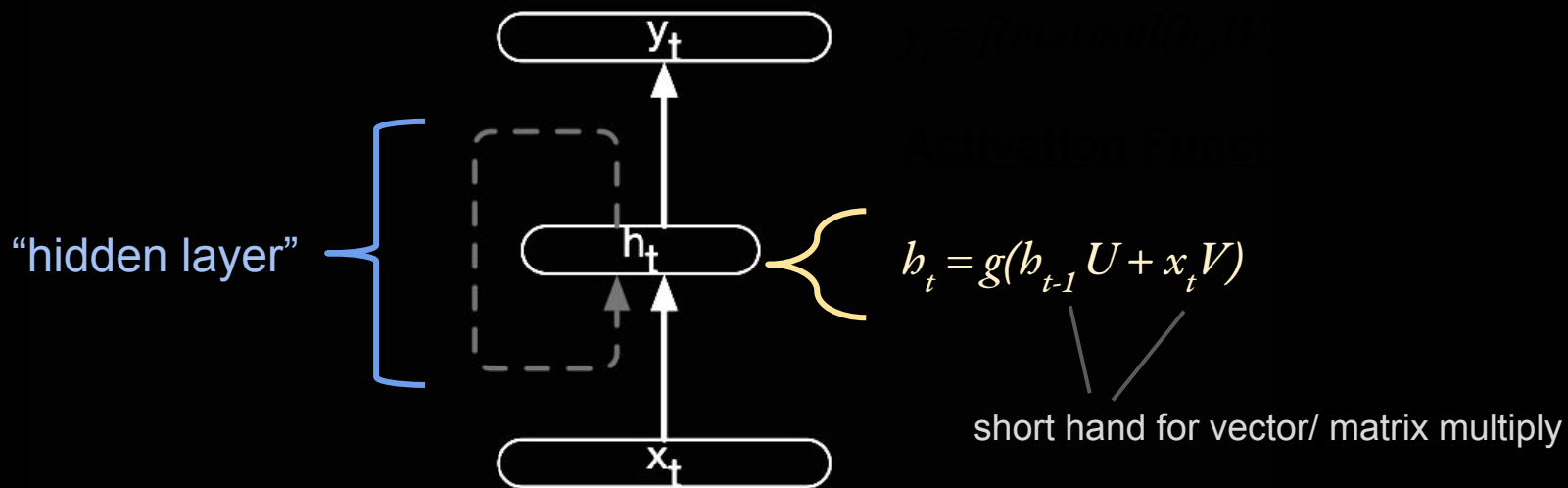


Figure 9.2 Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

Neural Networks: Graphs of Operations (excluding the optimization nodes)

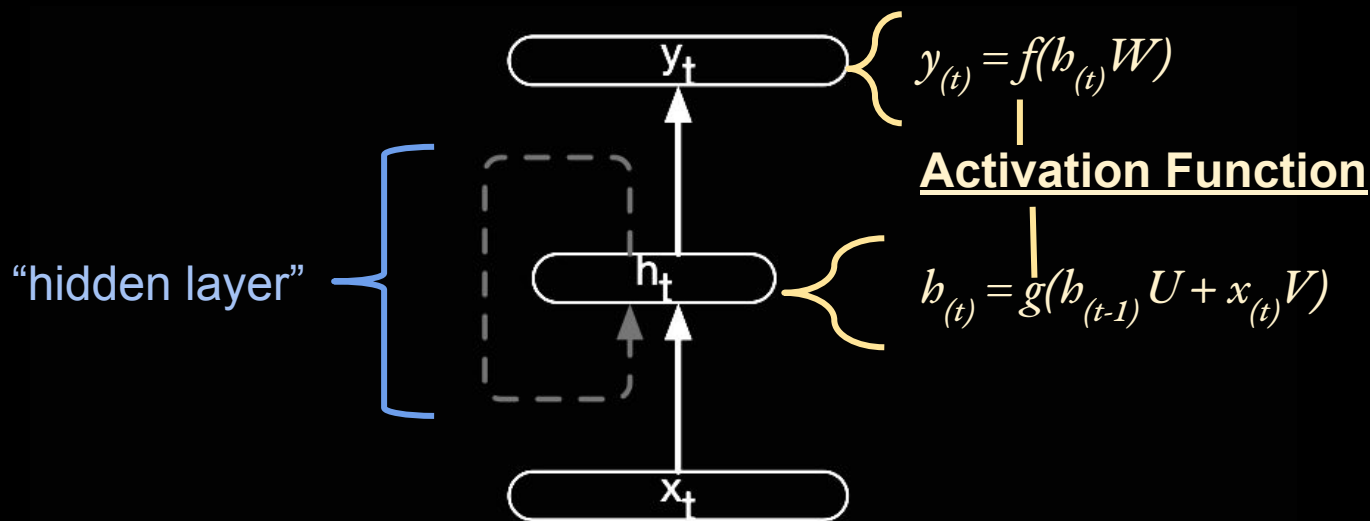
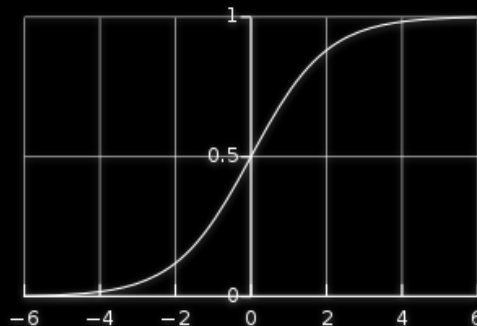


Figure 9.2 Simple recurrent neural network after Elman (Elman, 1990). The hidden layer includes a recurrent connection as part of its input. That is, the activation value of the hidden layer depends on the current input as well as the activation value of the hidden layer from the previous timestep. (Jurafsky, 2019)

Common Activation Functions

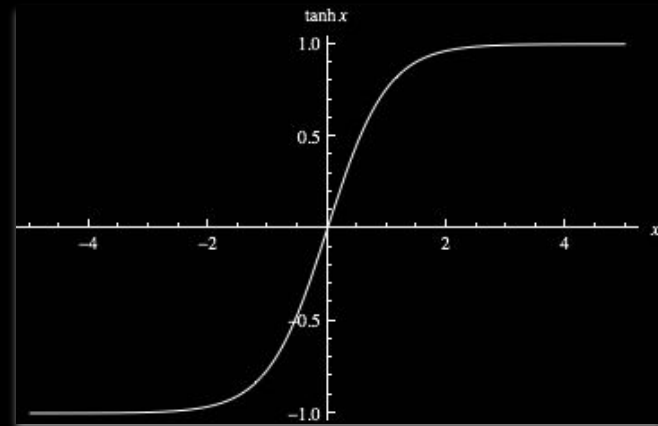
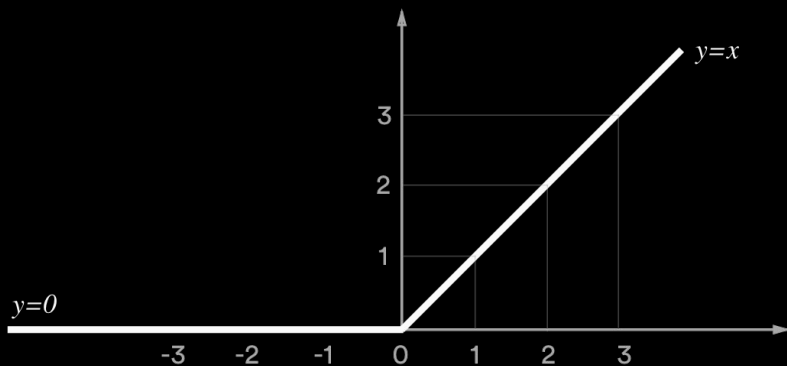
$$z = b_{(t)}W$$

Logistic: $\sigma(z) = 1 / (1 + e^{-z})$



Hyperbolic tangent: $\tanh(z) = 2\sigma(2z) - 1 = (e^{2z} - 1) / (e^{2z} + 1)$

Rectified linear unit (ReLU): $ReLU(z) = \max(0, z)$



Example: Forward Pass



(Geron, 2017)

```
#define forward pass graph:
```

```
 $h_{(0)} = \theta$ 
```

```
for  $i$  in range(1, len(x)):
```

```
     $h_{(i)} = g(U h_{(i-1)} + W x_{(i)})$  #update hidden state
```

```
     $y_{(i)} = f(V h_{(i)})$  #update output
```

Example: Forward Pass



```
#define forward pass graph:
```

```
 $h_{(0)} = \theta$ 
```

```
for  $i$  in range(1, len(x)):
```

```
     $h_{(i)} = g(U h_{(i-1)} + W x_{(i)})$  #update hidden state
```

```
     $y_{(i)} = f(V h_{(i)})$  #update output
```

Example: Forward Pass



```
#define forward pass graph:
```

```
 $h_{(0)} = \theta$ 
```

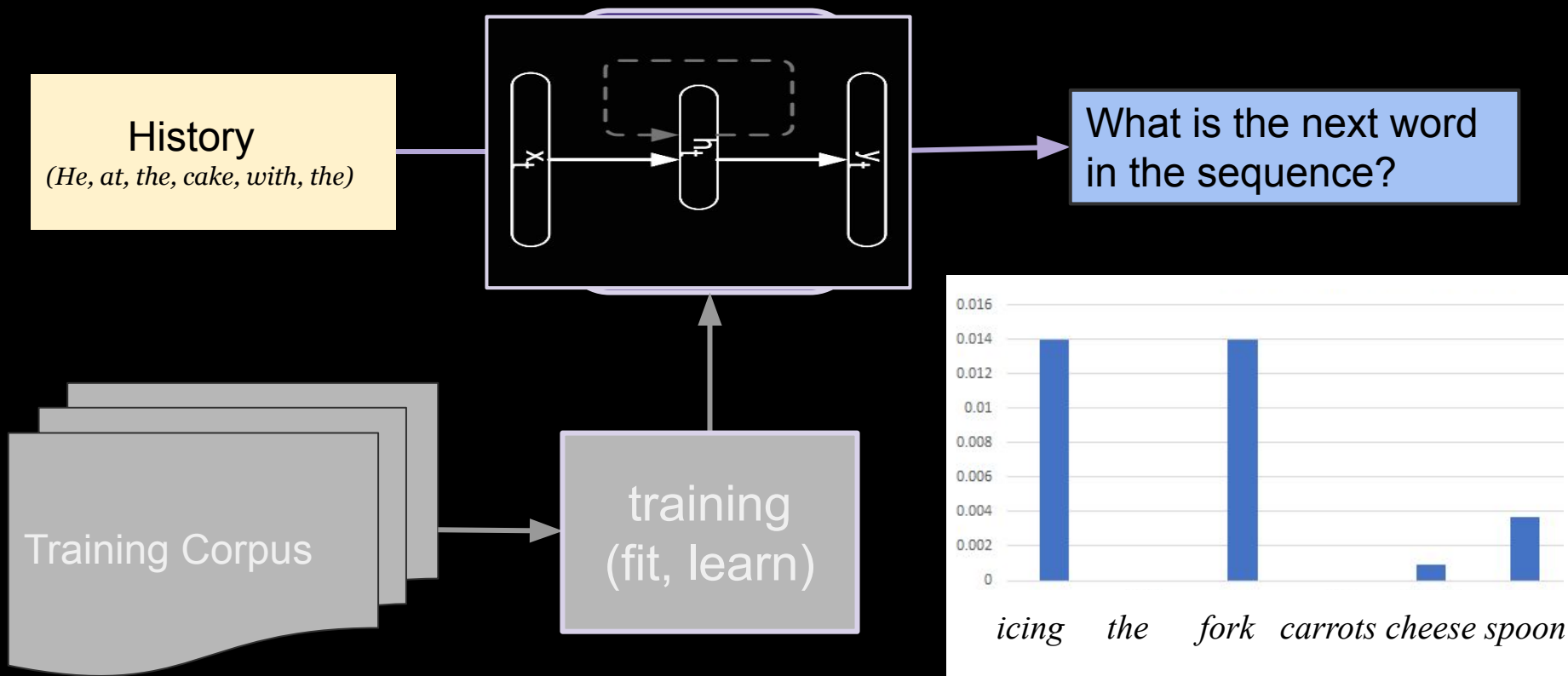
```
for  $i$  in range(1, len(x)):
```

```
     $h_{(i)} = \tanh(\text{matmul}(U, h_{(i-1)}) + \text{matmul}(W, x_{(i)}))$  #update hidden state
```

```
     $y_{(i)} = \text{softmax}(\text{matmul}(V, h_{(i)}))$  #update output
```

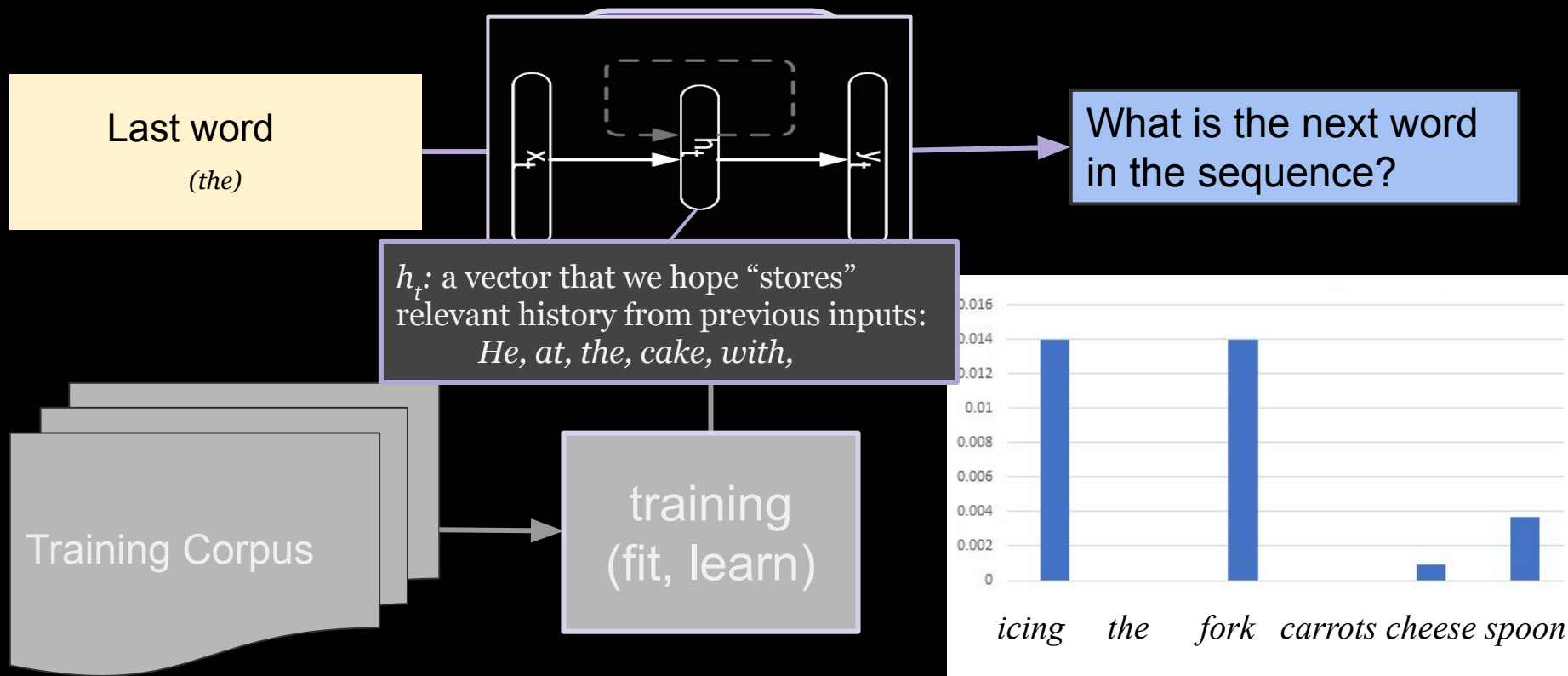
Language Modeling

Task: Estimate $P(w_n | w_1, w_2, \dots, w_{n-1})$
:probability of a next word given history
 $P(\text{fork} | \text{He ate the cake with the}) = ?$



Language Modeling

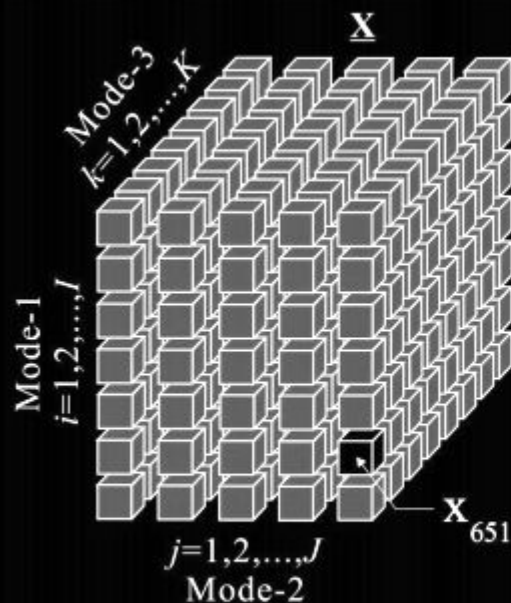
Task: Estimate $P(w_n | w_1, w_2, \dots, w_{n-1})$
:probability of a next word given history
 $P(\text{fork} | \text{He ate the cake with the}) = ?$



Tensors in PyTorch

Need a workflow system catered to numerical computation.

Basic idea: defines a graph of operations on tensors



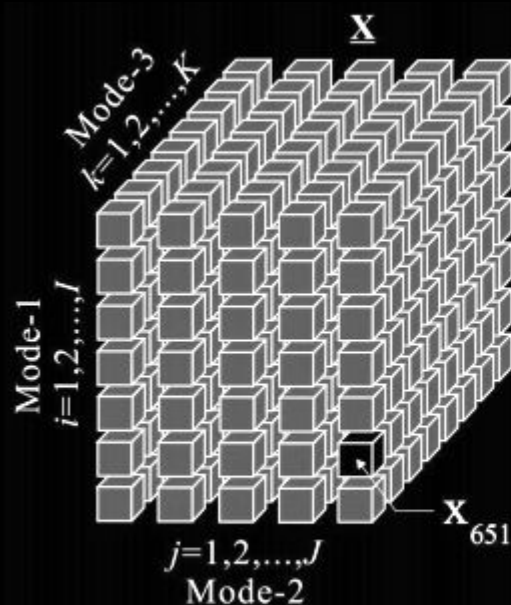
Tensors

Need a workflow system catered to numerical computation.

Basic idea: defines a graph of operations on tensors



A multi-dimensional matrix



PyTorch

A workflow system catered to numerical computation.

Basic idea: defines a graph of operations on tensors

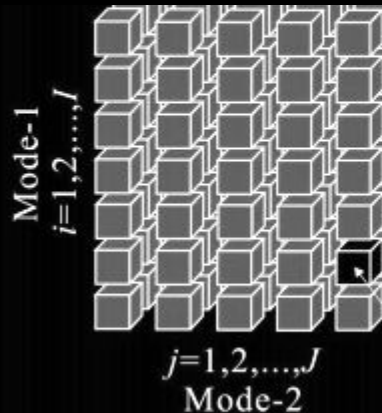


A multi-dimensional matrix

A 2-d tensor is just a matrix.

1-d: vector

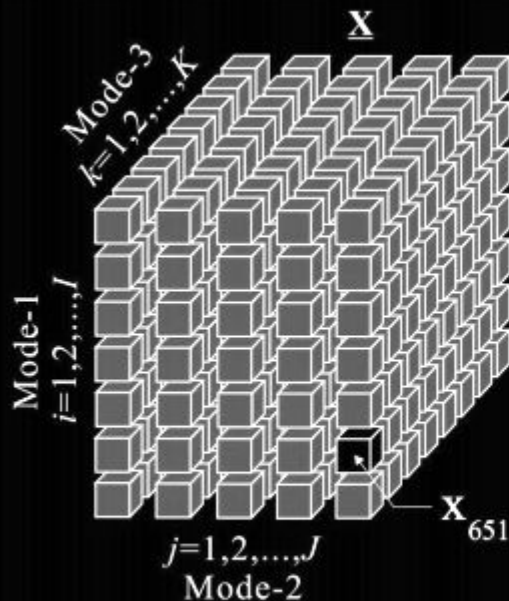
0-d: a constant / scalar



PyTorch

A workflow system catered to numerical computation.

Basic idea: defines a graph of operations on tensors



A multi-dimensional matrix

A 2-d tensor is just a matrix.

1-d: vector

0-d: a constant / scalar

Linguistic Ambiguity:

“ds” of a Tensor \neq

Dimensions of a Matrix

(i.stack.imgur.com)

PyTorch

A workflow system catered to numerical computation.

Basic idea: defines a graph of operations on **tensors**

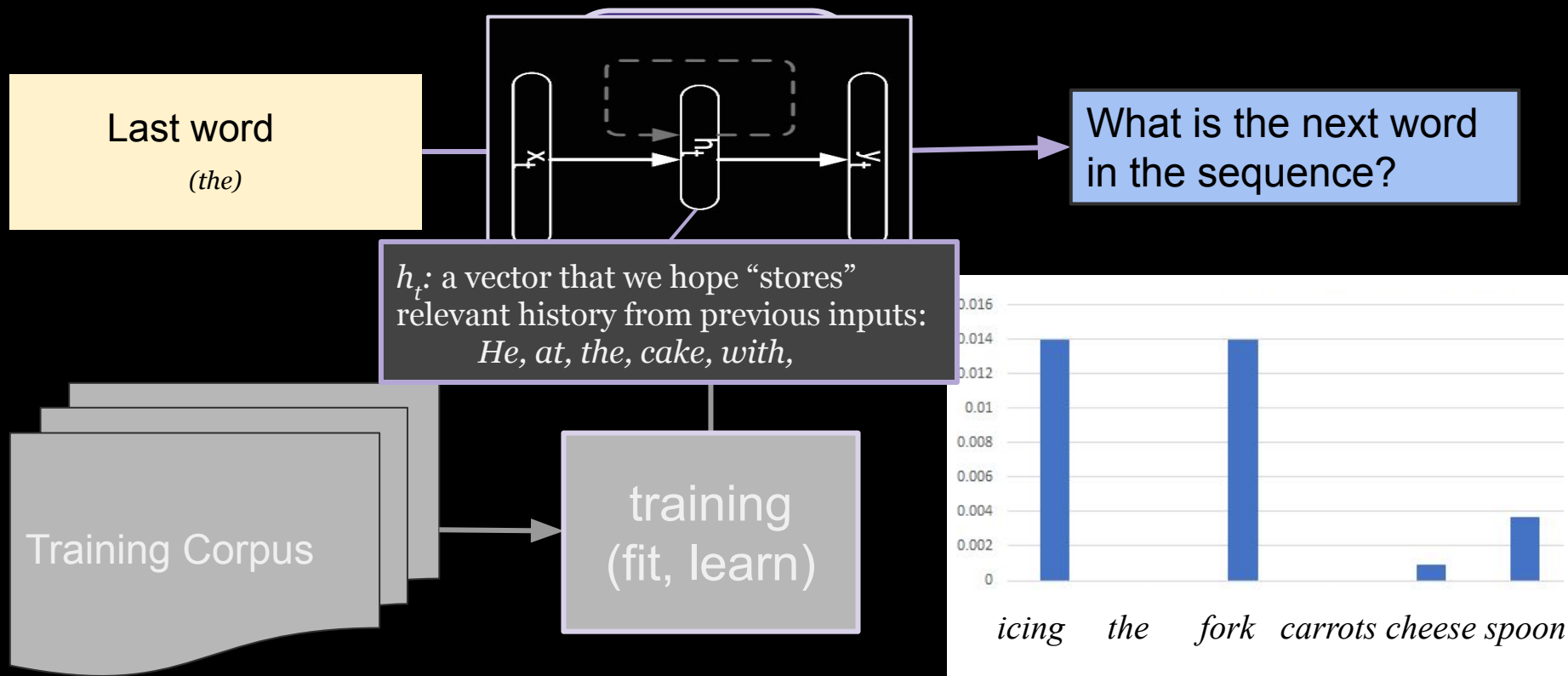
Why?

Efficient, high-level built-in **linear algebra** and **machine learning optimization operations** (i.e. transformations).

enables complex models, like deep learning

Language Modeling

Task: Estimate $P(w_n | w_1, w_2, \dots, w_{n-1})$
:probability of a next word given history
 $P(\text{fork} | \text{He ate the cake with the}) = ?$



Example: RNN



```
def forward(self, X):
    #Basic RNN Forward Pass:
     $h_{(0)} = 0$ 
    for i in range(1, len(x)):
         $h_{(i)} = \text{torch.tanh}(\text{torch.matmul}(U, h_{(i-1)}) + \text{torch.matmul}(W, x_{(i)}))$  #update
        hidden state
         $y_{(i)} = \text{nn.log_softmax}(\text{torch.matmul}(V, h_{(i)}))$  #update output

    ...

loss_func = nn.NLLLoss() #normalized log likelihood loss
            #torch.mean(-torch.sum(y*y_pred))
```

Example: RNN



```
def forward(self, X):
```

```
    #Basic RNN Forward Pass:
```

y_t

$$y_{(t)} = f(h_{(t)} W)$$

Activation Function $u_1(W, x_{(i)})$ #update

h_t

$$h_{(t)} = g(h_{(t-1)} U + x_{(t)} V) \text{ output}$$

x_t

“hidden layer”

```
    loss_func = nn.NLLLoss() #negative log likelihood loss
```

```
    #torch.mean(-torch.sum(y*y_pred))
```

Example: RNN

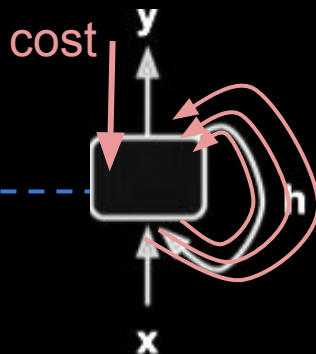


```
def forward(self, X):
    #Basic RNN Forward Pass:
     $h_{(0)} = 0$ 
    for i in range(1, len(x)):
         $h_{(i)} = \text{torch.tanh}(\text{torch.matmul}(U, h_{(i-1)}) + \text{torch.matmul}(W, x_{(i)}))$  #update
        hidden state
         $y_{(i)} = \text{nn.log_softmax}(\text{torch.matmul}(V, h_{(i)}))$  #update output

    ...

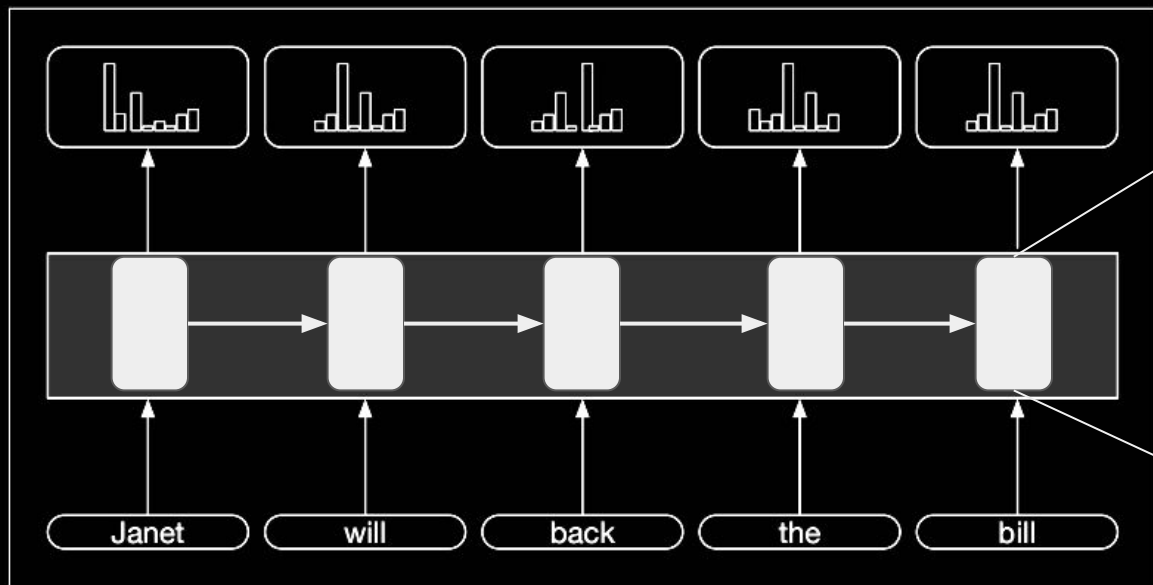
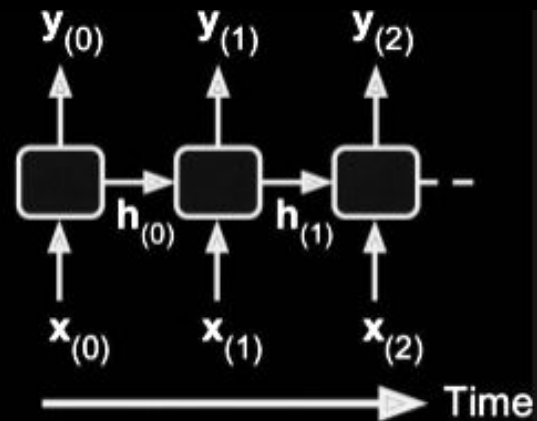
loss_func = nn.NLLLoss() #negative log likelihood loss
            #torch.mean(-torch.sum(y*y_pred))
```

Back Propagation



```
def forward(self, X):  
    #Basic RNN Forward Pass:  
     $h_{(0)} = 0$   
    for i in range(1, len(x)):  
         $h_{(i)} = \text{torch.tanh}(\text{torch.matmul}(U, h_{(i-1)}) + \text{torch.matmul}(W, x_{(i)}))$  #update  
        hidden state  
         $y_{(i)} = \text{nn.log_softmax}(\text{torch.matmul}(V, h_{(i)}))$  #update output  
  
    ...  
  
    loss_func = nn.NLLLoss() #negative log likelihood loss  
    #torch.mean(-torch.sum(y*y_pred))
```


Solution: Unrolling



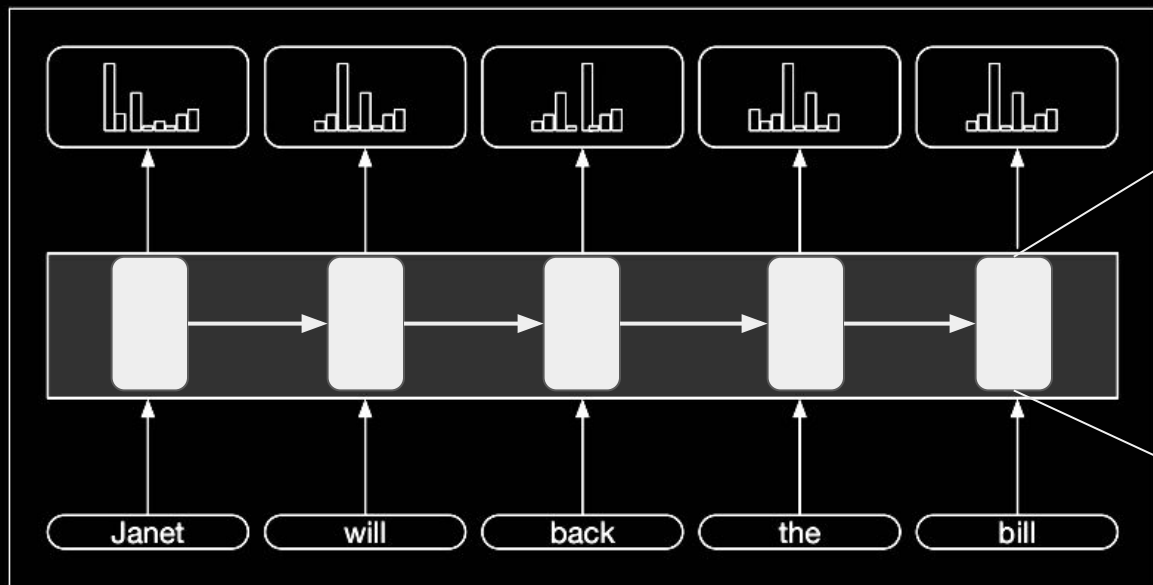
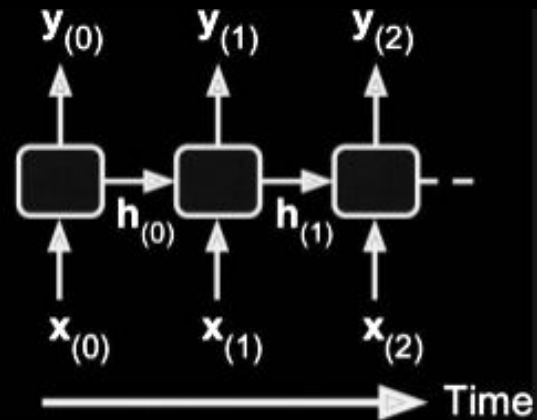
$$y_{(t)} = f(h_{(t)}W)$$

Activation Function

$$h_{(t)} = g(h_{(t-1)}U + x_{(t)}V)$$

Figure 9.8 Part-of-speech tagging as sequence labeling with a simple RNN. Pre-trained word embeddings serve as inputs and a softmax layer provides a probability distribution over the part-of-speech tags as output at each time step.

Solution: Unrolling



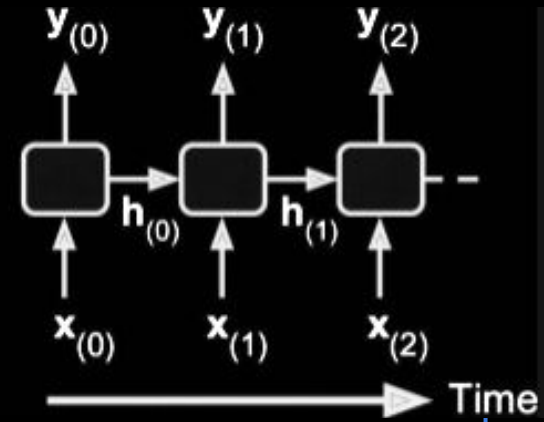
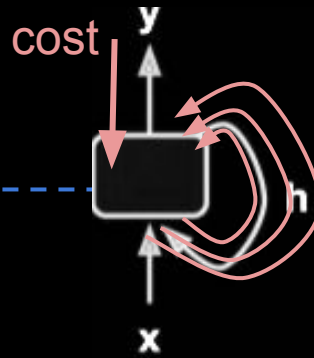
$$y_{("bill")} = f(h_{("bill")}W)$$

Activation Function

$$h_{("bill")} = g(h_{("the")}U + x_{("bill")}V)$$

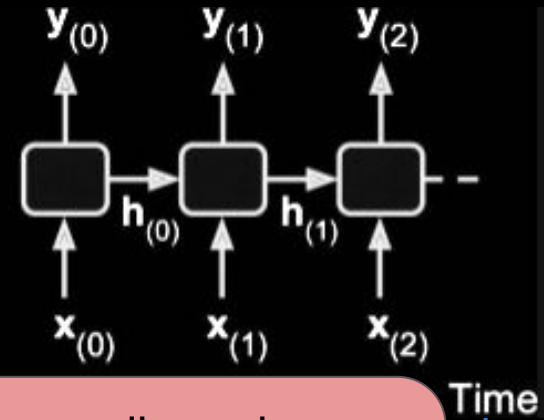
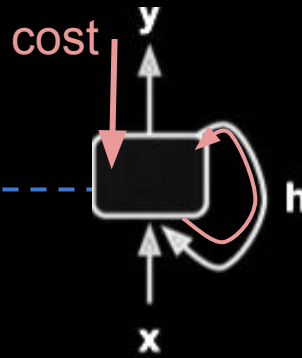
Figure 9.8 Part-of-speech tagging as sequence labeling with a simple RNN. Pre-trained word embeddings serve as inputs and a softmax layer provides a probability distribution over the part-of-speech tags as output at each time step.

Back Propagation



```
def forward(self, X):  
    #Basic RNN Forward Pass:  
     $h_{(0)} = 0$   
    for i in range(1, len(x)):  
         $h_{(i)} = \text{torch.tanh}(\text{torch.matmul}(U, h_{(i-1)}) + \text{torch.matmul}(W, x_{(i)}))$  #update  
        hidden state  
         $y_{(i)} = \text{nn.log_softmax}(\text{torch.matmul}(V, h_{(i)}))$  #update output  
  
    ...  
  
    loss_func = nn.NLLLoss() #negative log likelihood loss  
    #torch.mean(-torch.sum(y*y_pred))
```

Back Propagation



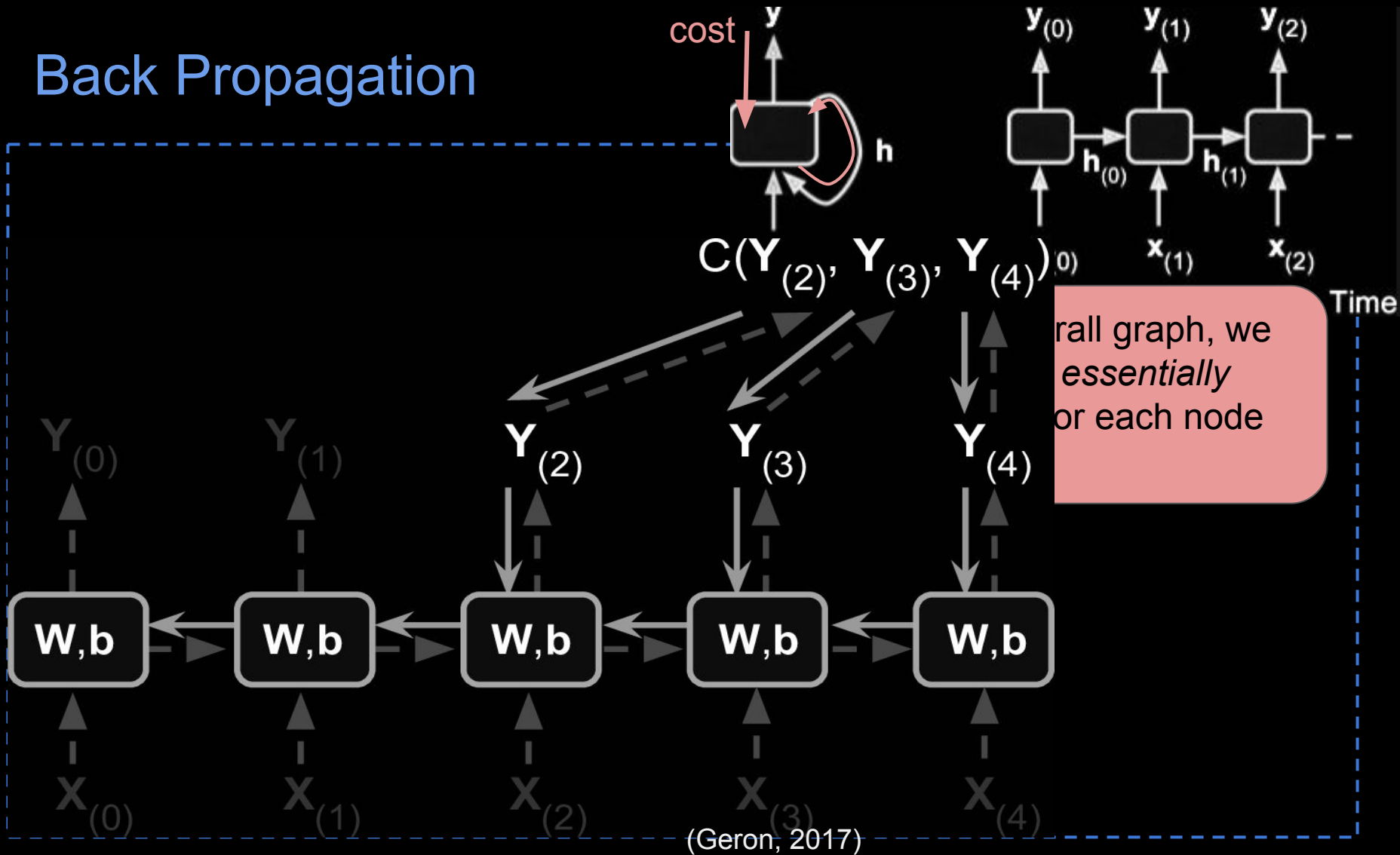
```
def forward(self, X):
    #Basic RNN Forward Pass:
    h(0) = 0
    for i in range(1, len(x)):
        h(i) = torch.tanh(torch.
            hidden state
        y(i) = nn.log_softmax(to
```

To find the gradient for the overall graph, we use **back propogation**, which *essentially* chains together the gradients for each node (function) in the graph.

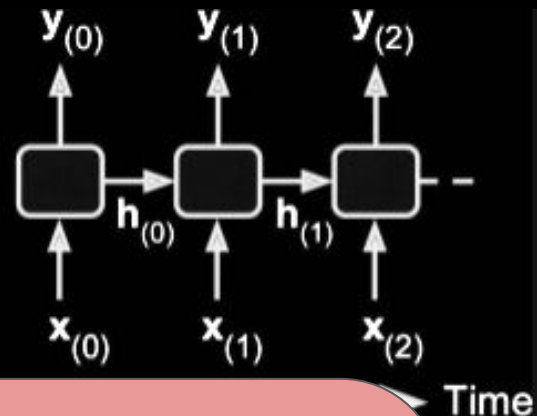
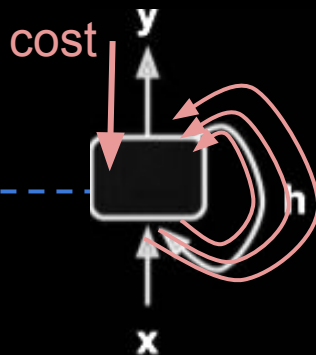
...

```
loss_func = nn.NLLLoss() #negative log likelihood loss
            #torch.mean(-torch.sum(y*y_pred))
```

Back Propagation



Back Propagation

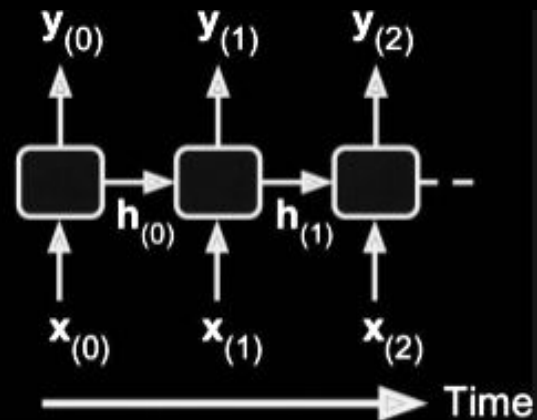
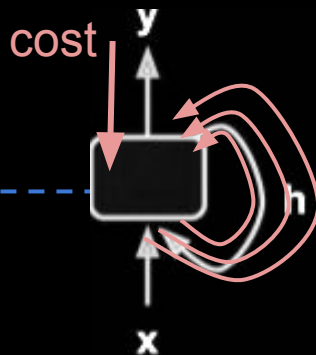


```
def forward(self, X):  
    #Basic RNN Forward Pass:  
     $h_{(0)} = 0$   
    for i in range(1, len(x)):  
         $h_{(i)} = \text{torch.tanh}(\text{torch.}$   
        hidden state  
         $y_{(i)} = \text{nn.log_softmax}(\text{to}$   
  
    ...  
  
    loss_func = nn.NLLLoss() #nega  
        #torch.mean(-torch.sum(y*y_pred))
```

To find the gradient for the overall graph, we use **back propogation**, which *essentially* chains together the gradients for each node (function) in the graph.

With many recursions, the gradients can vanish or explode (become too large or small for floating point operations).

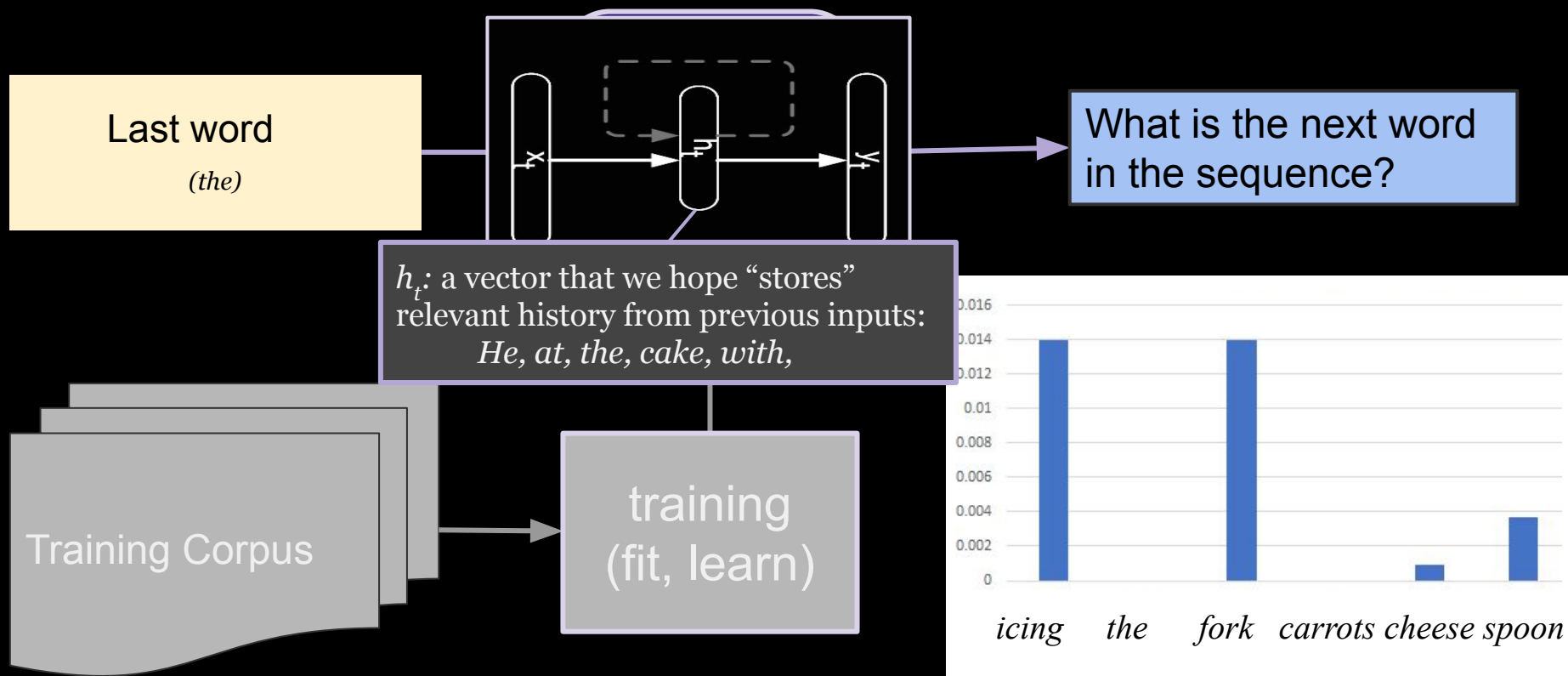
Back Propagation



```
def forward(self, X):  
    #Basic RNN Forward Pass:  
     $h_{(0)} = 0$   
    for i in range(1, len(x)):  
         $h_{(i)} = \text{torch.tanh}(\text{torch.matmul}(U, h_{(i-1)}) + \text{torch.matmul}(W, x_{(i)}))$  #update  
        hidden state  
         $y_{(i)} = \text{nn.log_softmax}(\text{torch.matmul}(V, h_{(i)}))$  #update output  
  
    ...  
  
loss_func = nn.NLLLoss() #negative log likelihood loss  
            #torch.mean(-torch.sum(y*y_pred))
```

How to address exploding and vanishing gradients?

How to address exploding and vanishing gradients?



How to address exploding and vanishing gradients?

Ad Hoc approaches: e.g. stop backprop iterations very early. “clip” gradients when too high.



How to address exploding and vanishing gradients?

Dominant approach: Use Long Short Term Memory Networks (LSTM)

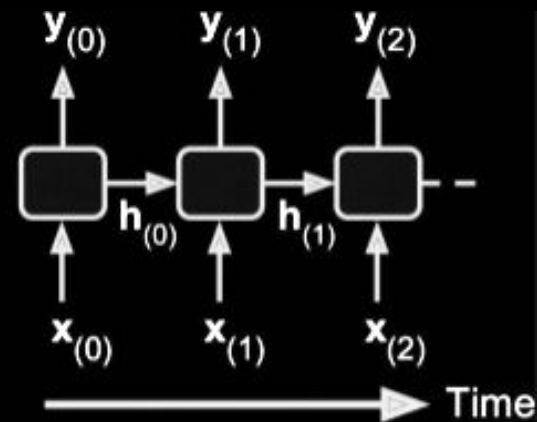


How to address exploding and vanishing gradients?

Dominant approach: Use Long Short Term Memory Networks (LSTM)



RNN model

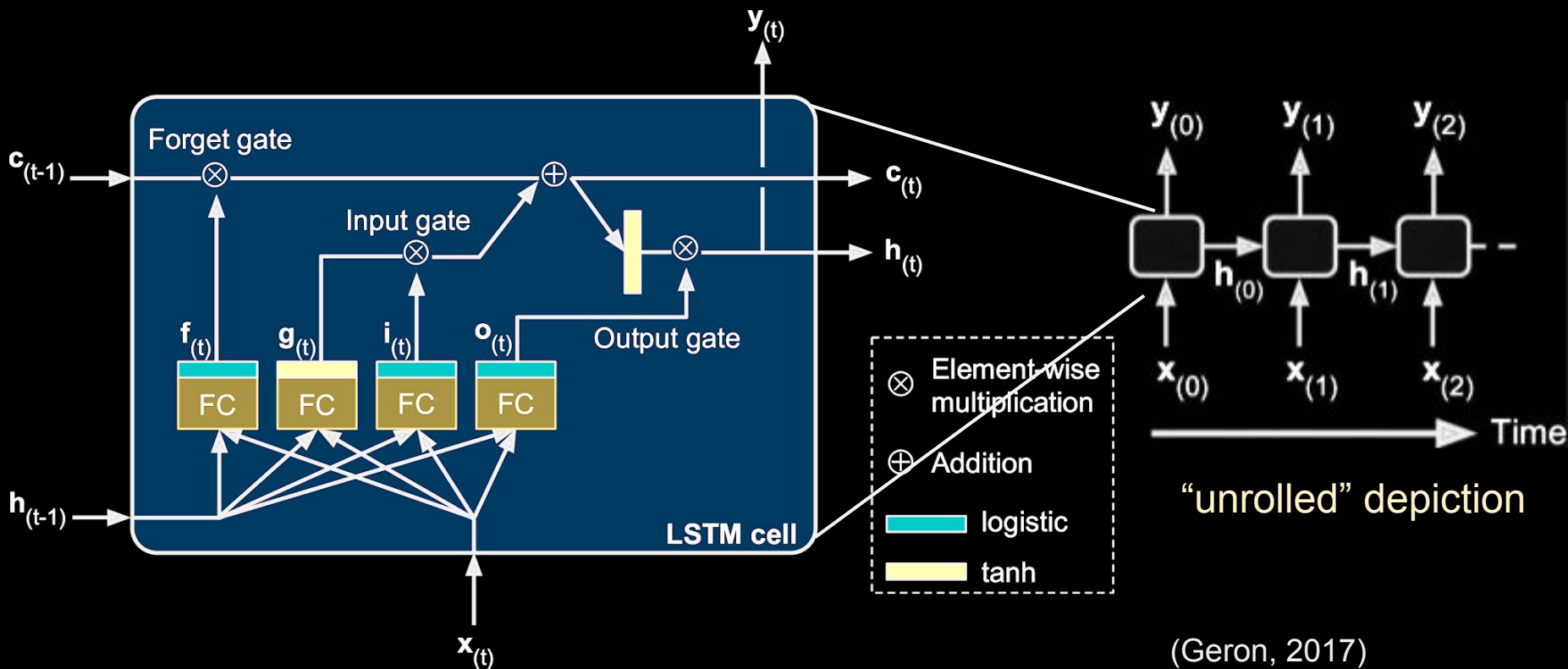


“unrolled” depiction

(Geron, 2017)

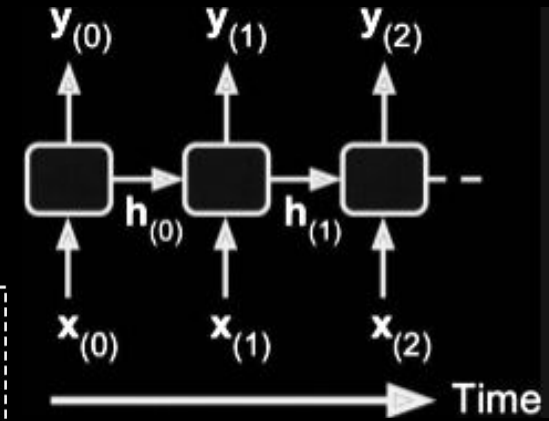
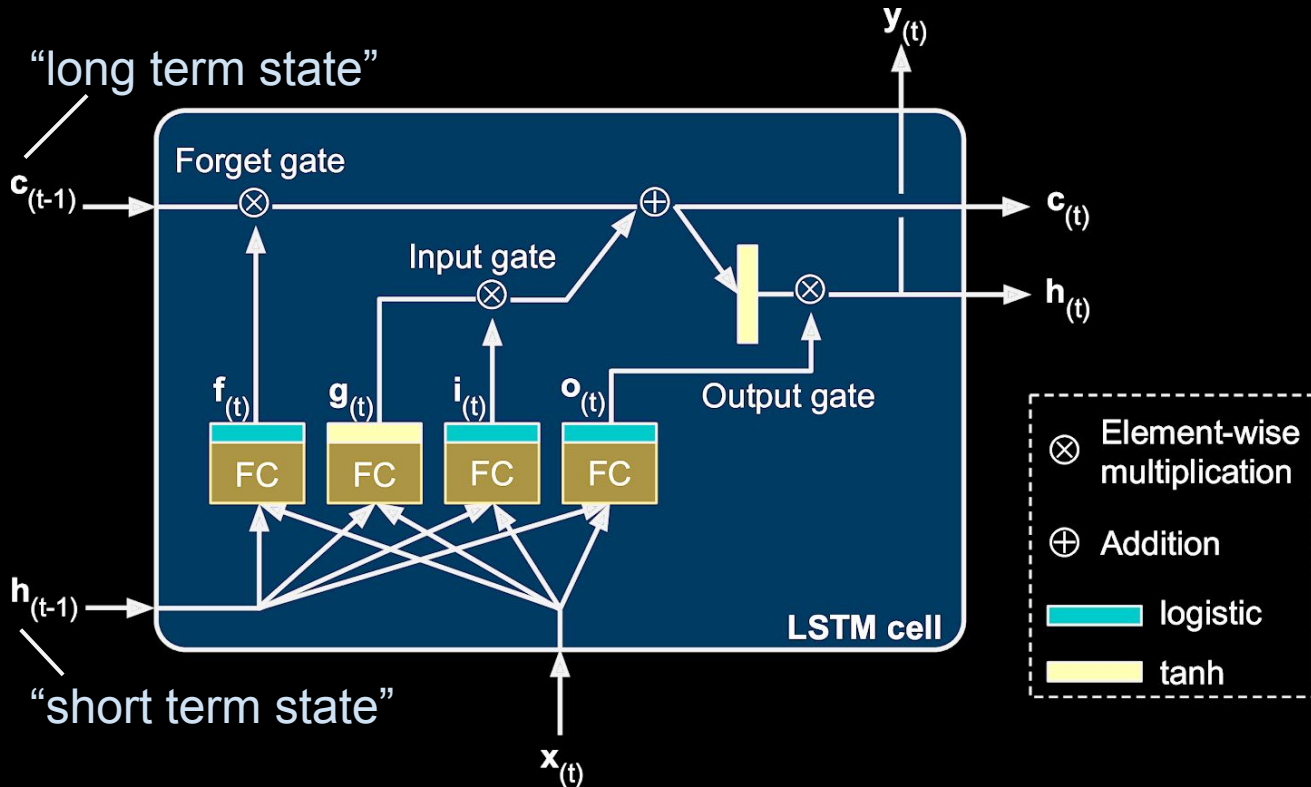
How to address exploding and vanishing gradients?

The LSTM Cell



How to address exploding and vanishing gradients?

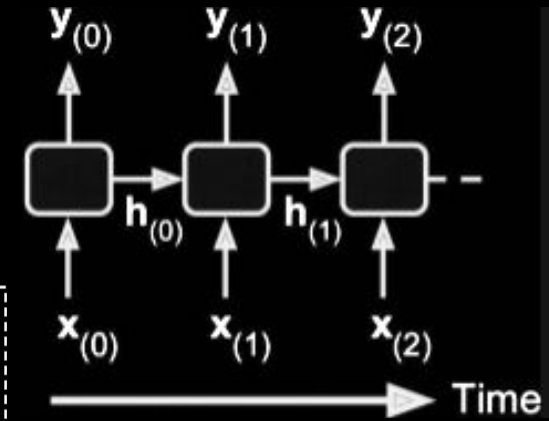
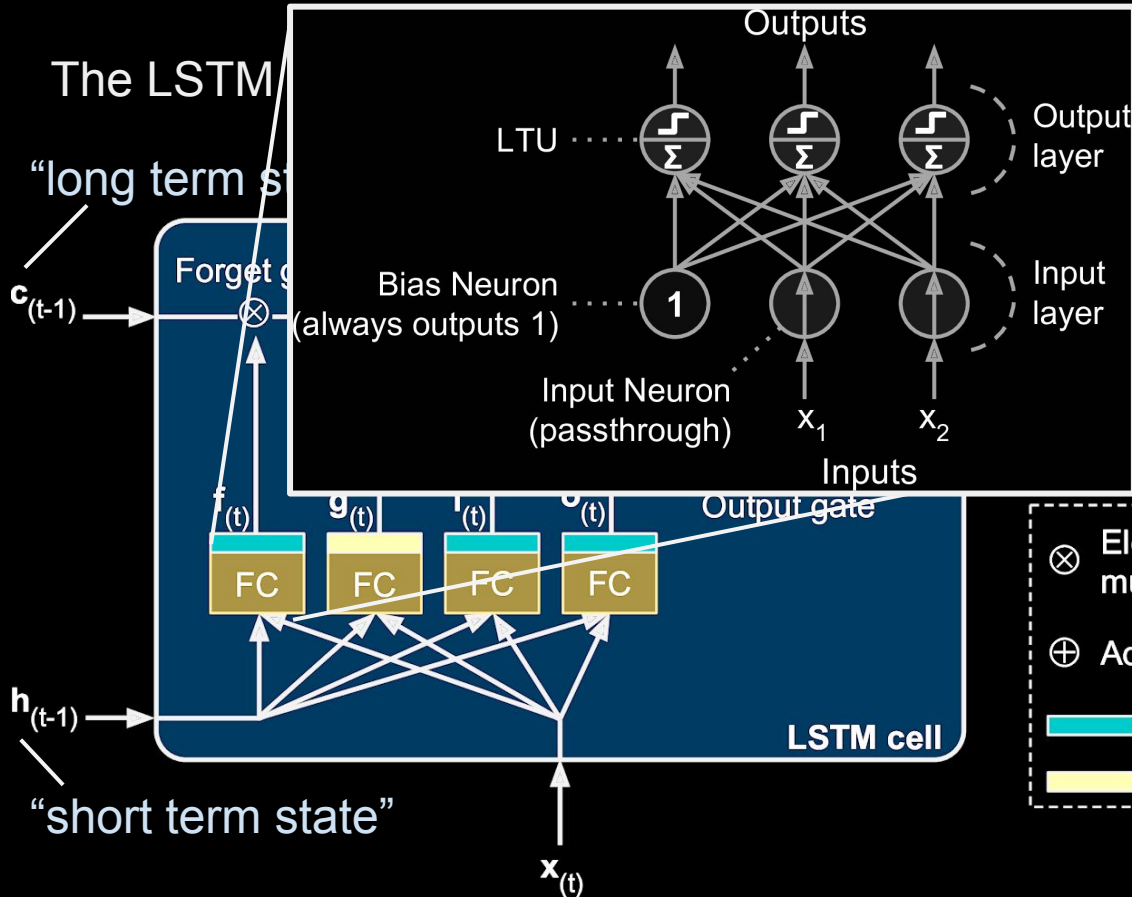
The LSTM Cell



"unrolled" depiction

(Geron, 2017)

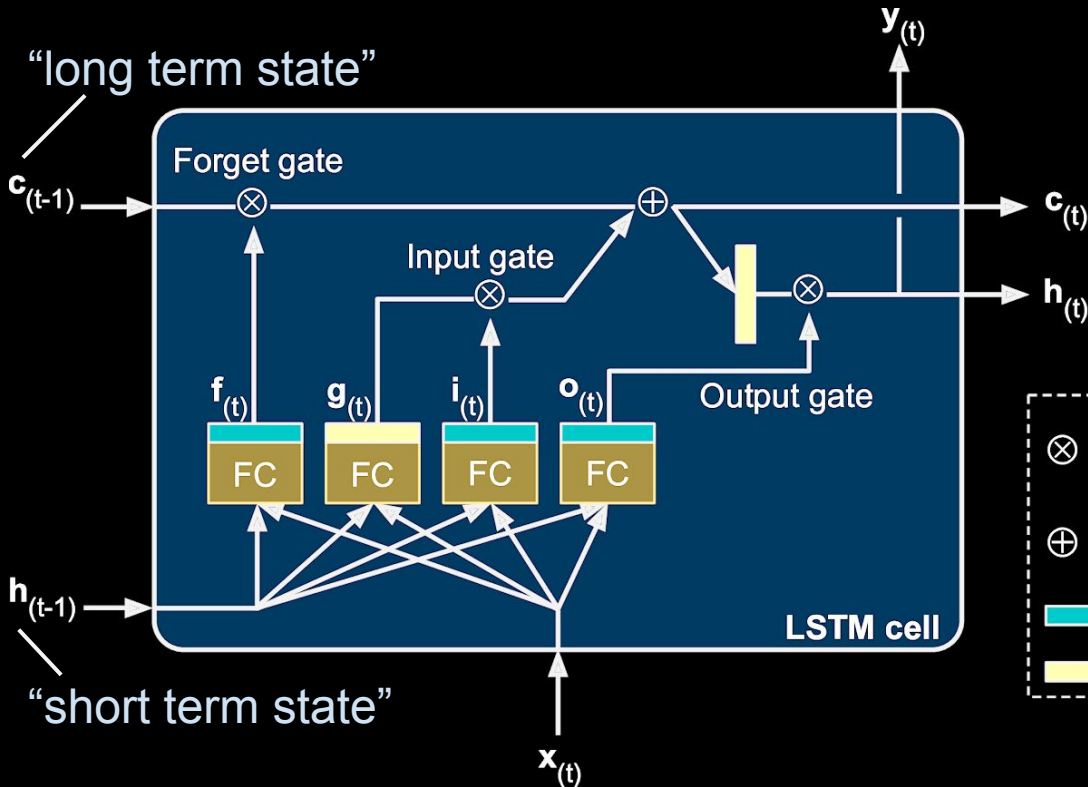
How to address exploding and vanishing gradients?



(Geron, 2017)

How to address exploding and vanishing gradients?

The LSTM Cell



$$\mathbf{i}_{(t)} = \sigma(\mathbf{W}_{xi}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hi}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{f}_{(t)} = \sigma(\mathbf{W}_{xf}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hf}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_f)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_g)$$

$$\mathbf{c}_{(t)} = \mathbf{f}_{(t)} \otimes \mathbf{c}_{(t-1)} + \mathbf{i}_{(t)} \otimes \mathbf{g}_{(t)}$$

bias term

\otimes Element-wise multiplication

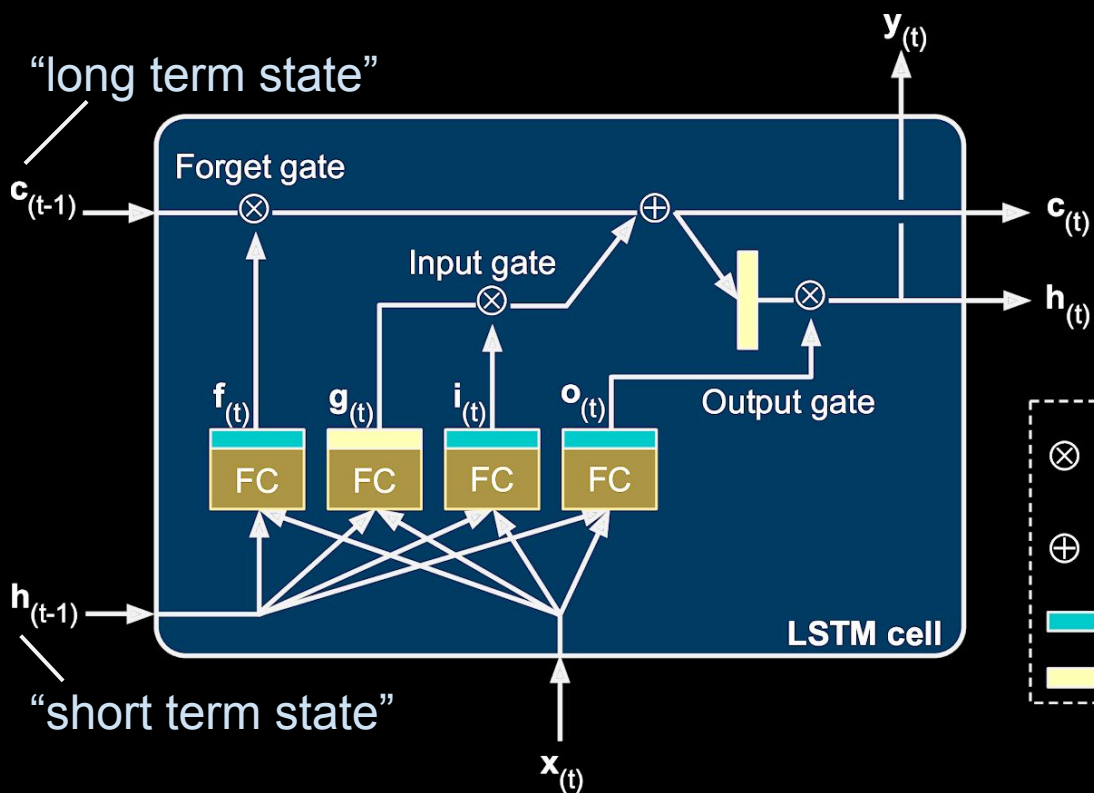
\oplus Addition

 logistic

 tanh

LSTM

The LSTM Cell



$$\mathbf{i}^{(t)} = \sigma(\mathbf{W}_{xi}^T \cdot \mathbf{x}^{(t)} + \mathbf{W}_{hi}^T \cdot \mathbf{h}^{(t-1)} + \mathbf{b}_i)$$

$$\mathbf{f}^{(t)} = \sigma(\mathbf{W}_{xf}^T \cdot \mathbf{x}^{(t)} + \mathbf{W}_{hf}^T \cdot \mathbf{h}^{(t-1)} + \mathbf{b}_f)$$

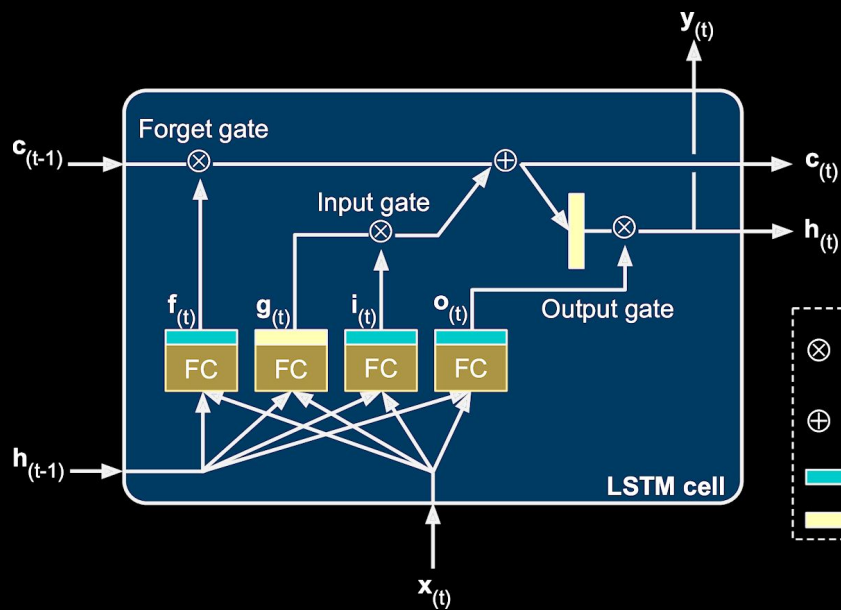
$$\mathbf{o}^{(t)} = \sigma(\mathbf{W}_{xo}^T \cdot \mathbf{x}^{(t)} + \mathbf{W}_{ho}^T \cdot \mathbf{h}^{(t-1)} + \mathbf{b}_o)$$

$$\mathbf{g}^{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}^{(t)} + \mathbf{W}_{hg}^T \cdot \mathbf{h}^{(t-1)} + \mathbf{b}_g)$$

$$\mathbf{c}^{(t)} = \mathbf{f}^{(t)} \otimes \mathbf{c}^{(t-1)} + \mathbf{i}^{(t)} \otimes \mathbf{g}^{(t)}$$

$$\mathbf{y}^{(t)} = \mathbf{h}^{(t)} = \mathbf{o}^{(t)} \otimes \tanh(\mathbf{c}^{(t)})$$

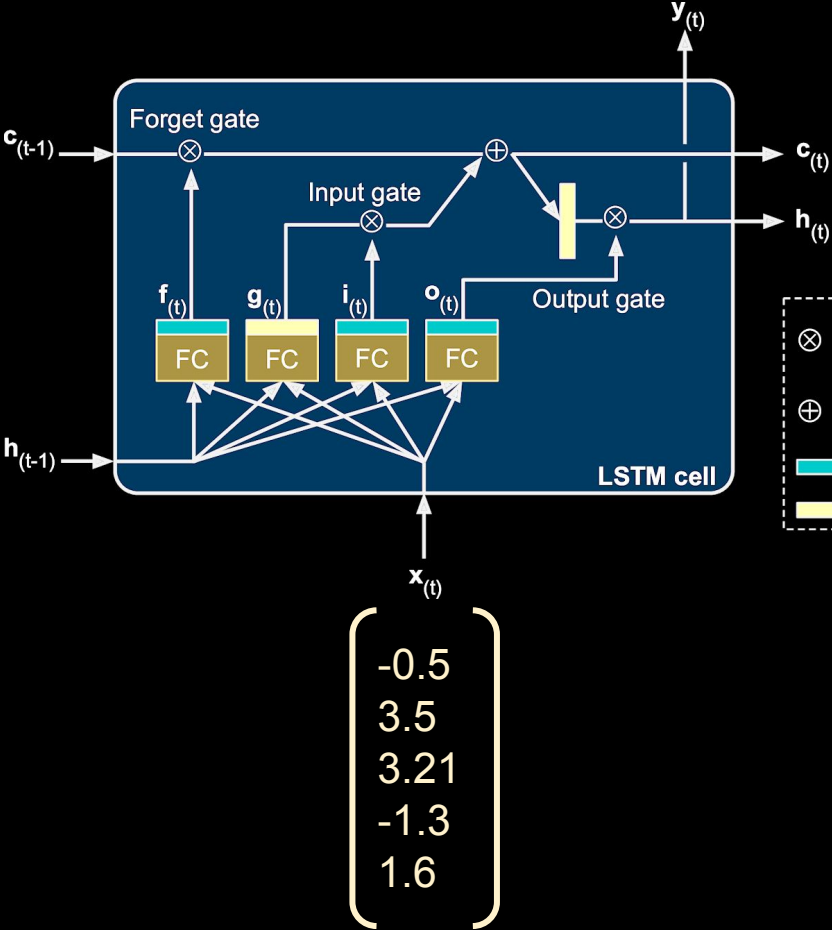
Input to LSTM



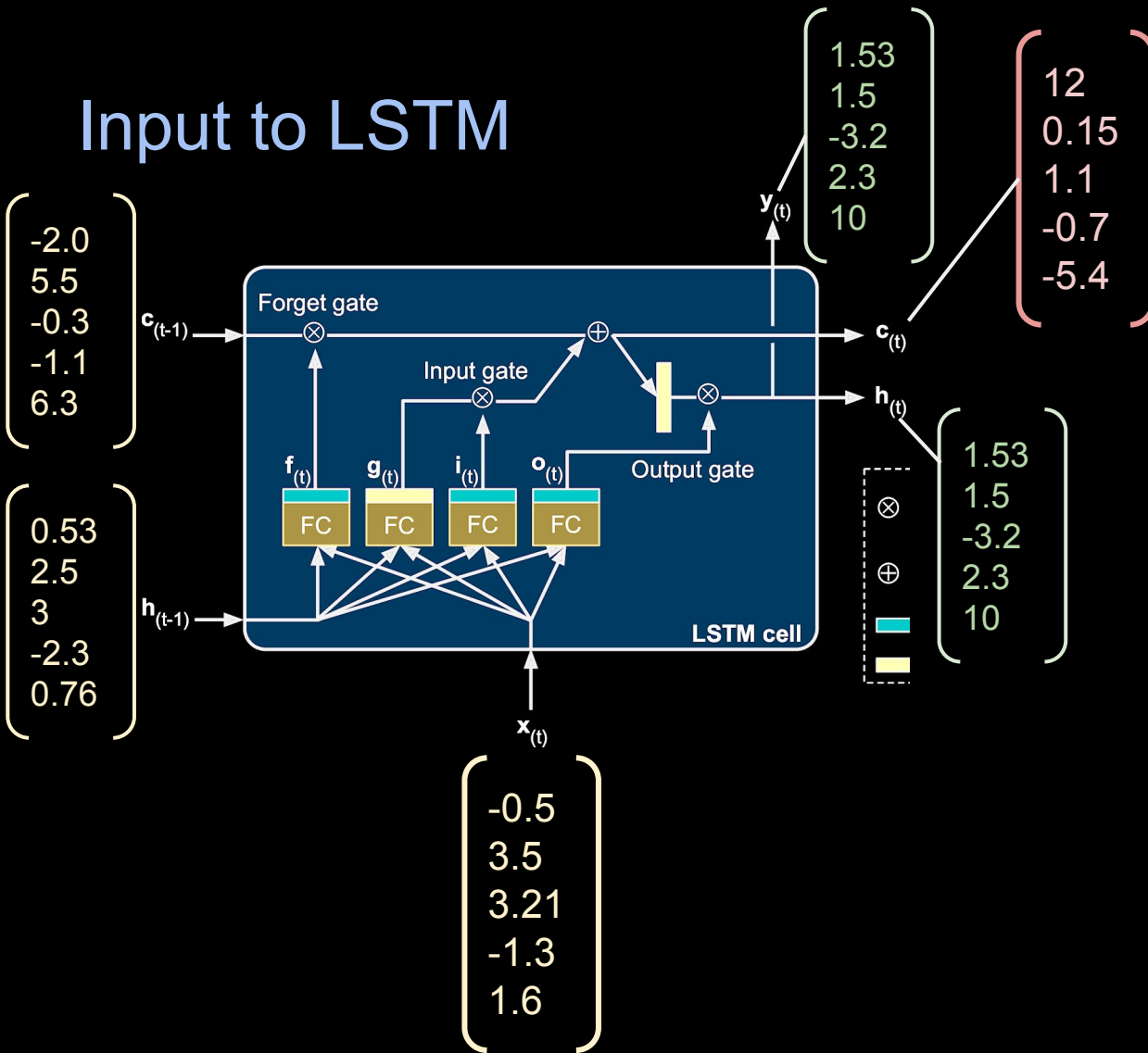
?

- One-hot encoding?
- Word Embedding

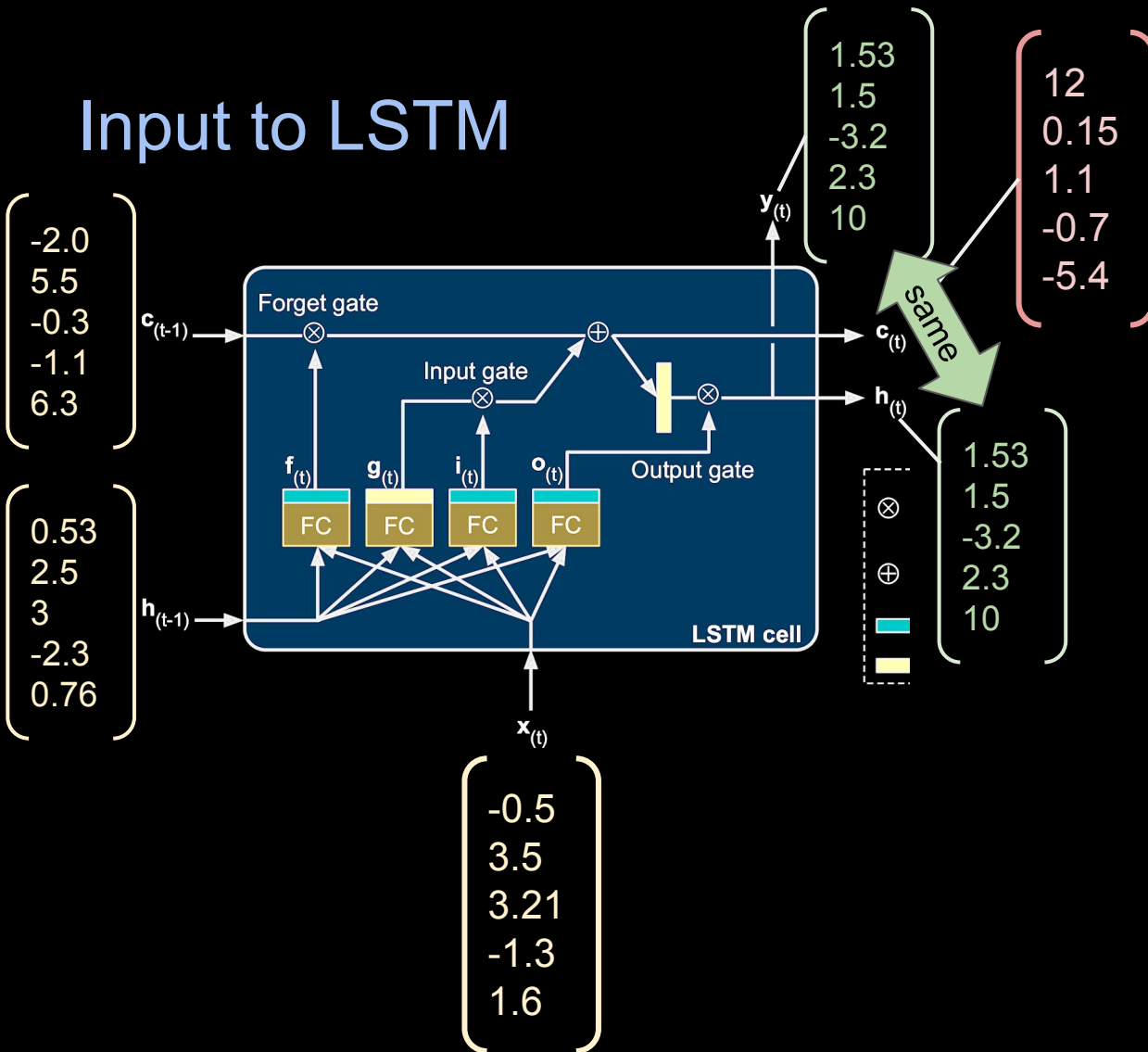
Input to LSTM



Input to LSTM

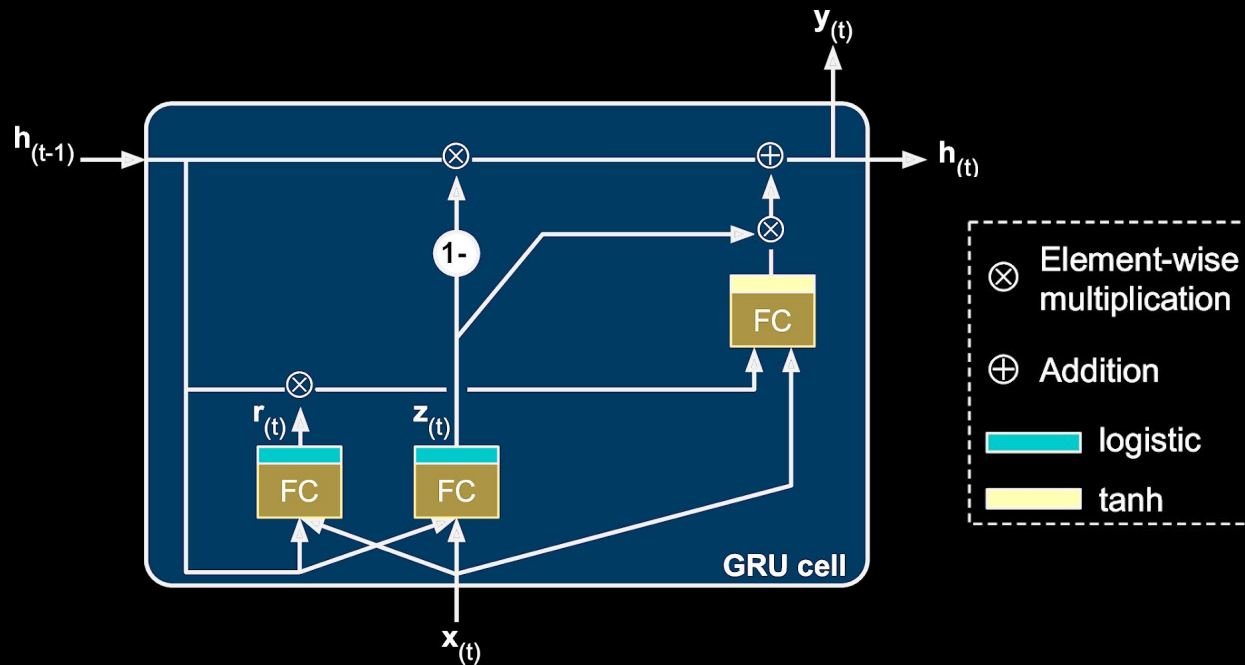


Input to LSTM



The GRU

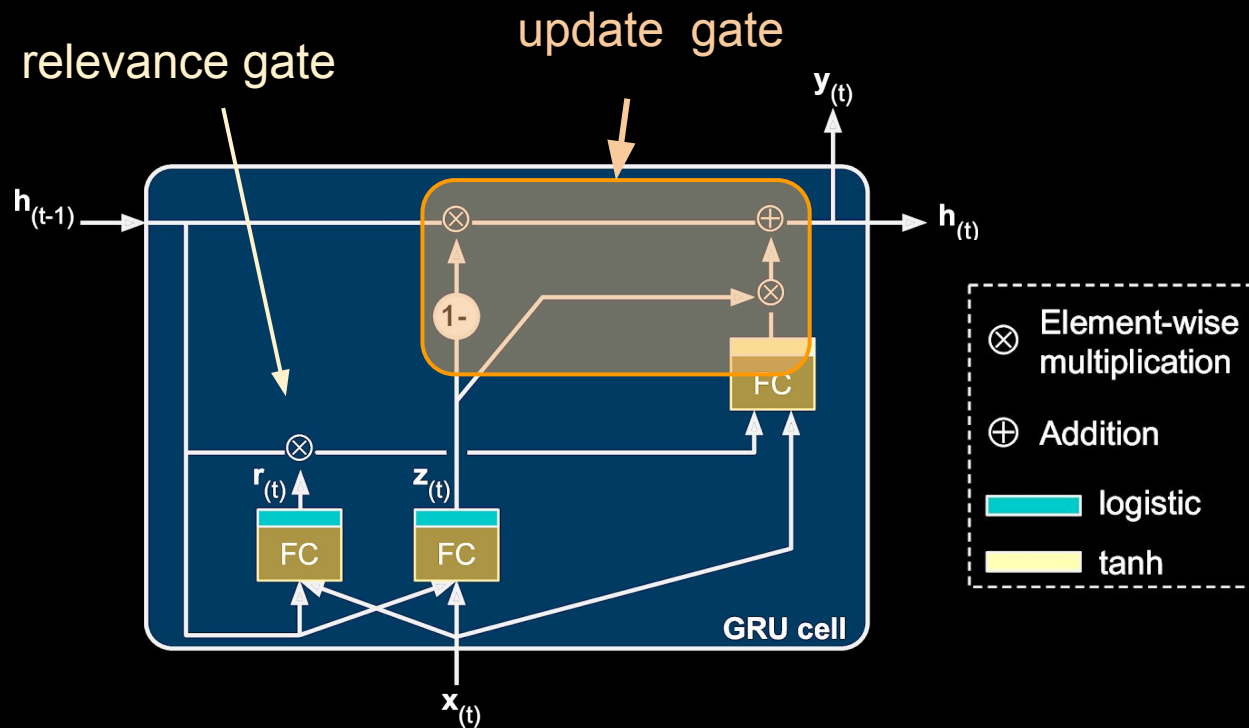
Gated Recurrent Unit



(Geron, 2017)

The GRU

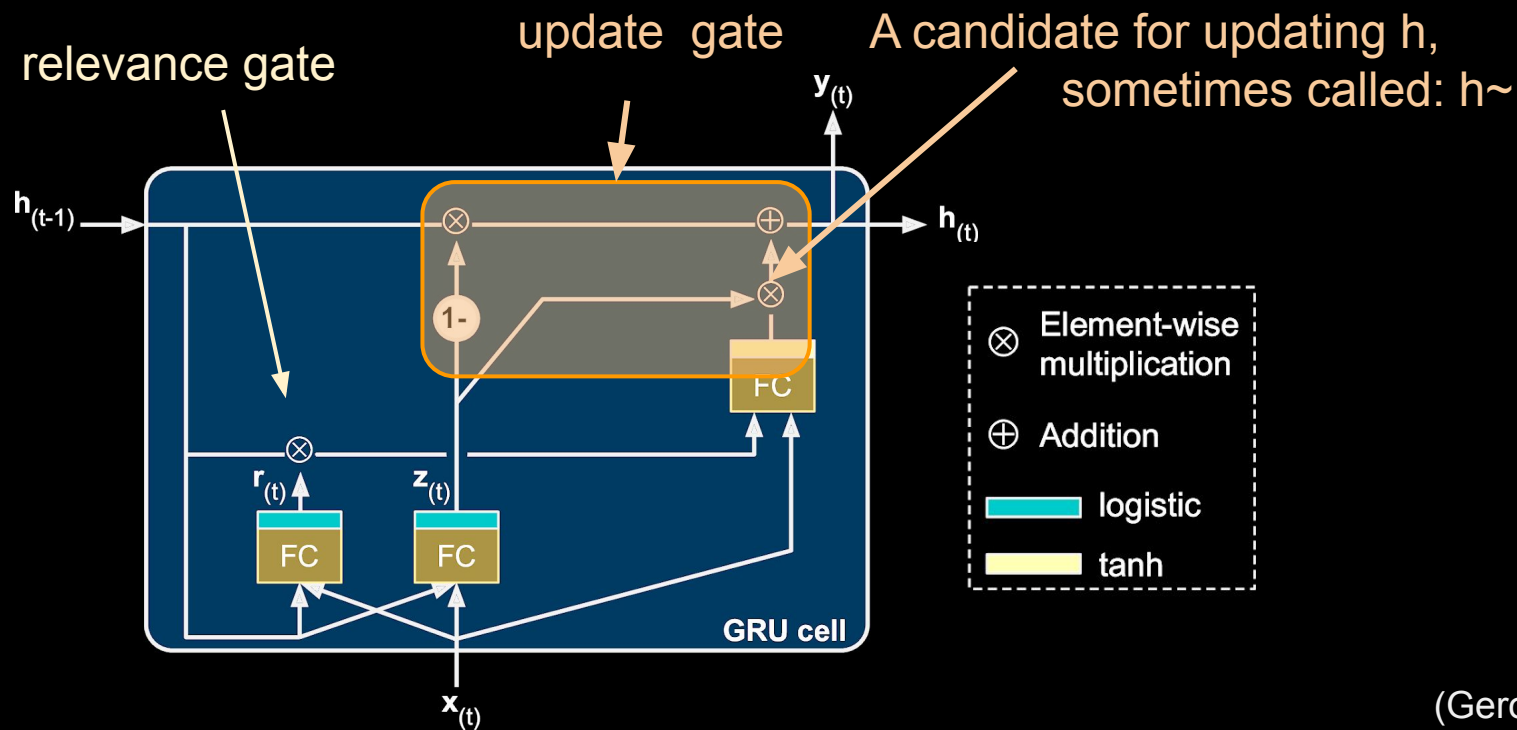
Gated Recurrent Unit



(Geron, 2017)

The GRU

Gated Recurrent Unit



(Geron, 2017)

The GRU

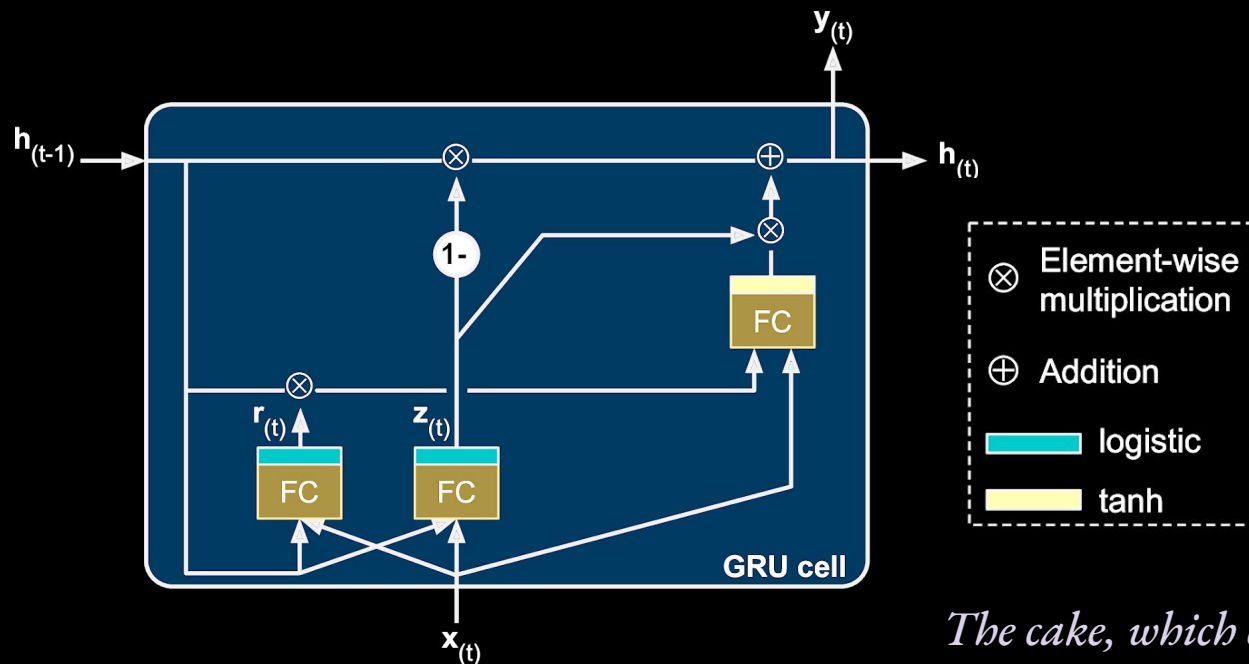
Gated Recurrent Unit

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$



The cake, which contained candles, was eaten.

What about the gradient?

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

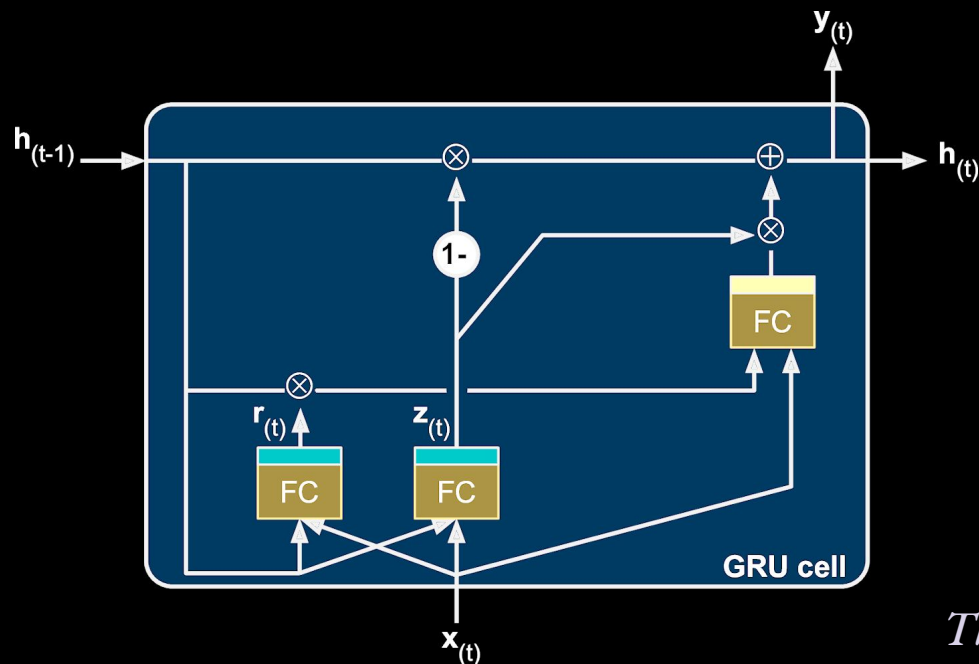
$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$

The gates (i.e. multiplications based on a logistic) often end up keeping the hidden state exactly (or nearly exactly) as it was. Thus, for most dimensions of \mathbf{h} ,

$$\mathbf{h}_{(t)} \approx \mathbf{h}_{(t-1)}$$



The cake, which contained candles, was eaten.

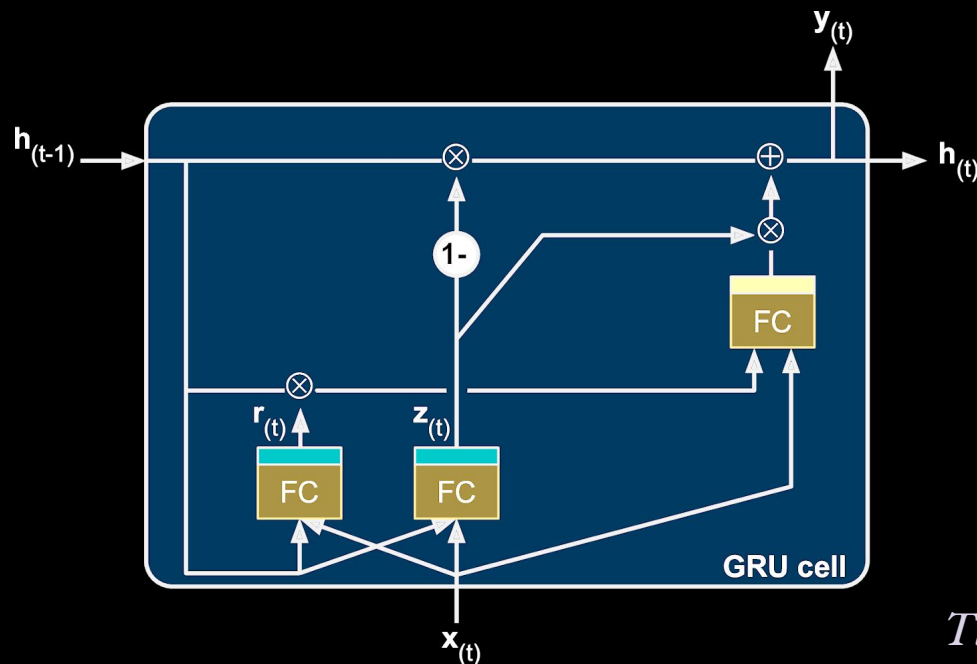
What about the gradient?

$$\mathbf{z}_{(t)} = \sigma(\mathbf{W}_{xz}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hz}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_z)$$

$$\mathbf{r}_{(t)} = \sigma(\mathbf{W}_{xr}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hr}^T \cdot \mathbf{h}_{(t-1)} + \mathbf{b}_r)$$

$$\mathbf{g}_{(t)} = \tanh(\mathbf{W}_{xg}^T \cdot \mathbf{x}_{(t)} + \mathbf{W}_{hg}^T \cdot (\mathbf{r}_{(t)} \otimes \mathbf{h}_{(t-1)}) + \mathbf{b}_g)$$

$$\mathbf{h}_{(t)} = \mathbf{z}_{(t)} \otimes \mathbf{h}_{(t-1)} + (1 - \mathbf{z}_{(t)}) \otimes \mathbf{g}_{(t)}$$



The gates (i.e. multiplications based on a logistic) often end up keeping the hidden state exactly (or nearly exactly) as it was. Thus, for most dimensions of \mathbf{h} ,

$$\mathbf{h}_{(t)} \approx \mathbf{h}_{(t-1)}$$

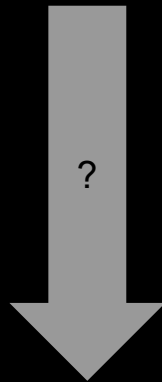
This tends to keep the gradient from vanishing since the same values will be present through multiple times in backpropagation through time. (The same idea applies to LSTMs but is easier to see here).

The cake, which contained candles, was eaten.

How to train an LSTM-style RNN

```
RNN_cost = torch.mean(-torch.sum(y*torch.log(y_pred)))  
#where did this come from?
```

Logistic Regression Likelihood: $L(\beta_0, \beta_1, \dots, \beta_k | X, Y) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$



Final Cost Function: $J^{(t)} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j}^{(t)} \log \hat{y}_{i,j}^{(t)}$ -- "cross entropy error"

How to train an LSTM-style RNN

```
RNN_cost = torch.mean(-torch.sum(y*torch.log(y_pred)))  
#where did this come from?
```

Logistic Regression Likelihood: $L(\beta_0, \beta_1, \dots, \beta_k | X, Y) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$

Log Likelihood: $\ell(\beta) = \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p(x_i))$

Final Cost Function: $J^{(t)} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j}^{(t)} \log \hat{y}_{i,j}^{(t)}$ -- "cross entropy error"

How to train an LSTM-style RNN

```
RNN_cost = torch.mean(-torch.sum(y*torch.log(y_pred)))  
#where did this come from?
```

Logistic Regression Likelihood: $L(\beta_0, \beta_1, \dots, \beta_k | X, Y) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$

Log Likelihood: $\ell(\beta) = \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p(x_i))$

Log Loss: $J(\beta) = -\frac{1}{N} \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p)(x_i)$

Final Cost Function: $J^{(t)} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j}^{(t)} \log \hat{y}_{i,j}^{(t)}$ -- "cross entropy error"

How to train an LSTM-style RNN

```
RNN_cost = torch.mean(-torch.sum(y*torch.log(y_pred)))  
#where did this come from?
```

Logistic Regression Likelihood: $L(\beta_0, \beta_1, \dots, \beta_k | X, Y) = \prod_{i=1}^n p(x_i)^{y_i} (1 - p(x_i))^{1-y_i}$

Log Likelihood: $\ell(\beta) = \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p(x_i))$

Log Loss: $J(\beta) = -\frac{1}{N} \sum_{i=1}^N y_i \log p(x_i) + (1 - y_i) \log (1 - p)(x_i)$

Cross-Entropy Cost: $J = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j} \log p(x_{i,j})$ (a "multiclass" log loss)

Final Cost Function: $J^{(t)} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|V|} y_{i,j}^{(t)} \log \hat{y}_{i,j}^{(t)}$ -- "cross entropy error"

How to train an LSTM-style RNN

```
RNN_cost = torch.mean(-torch.sum(y*torch.log(y_pred)))  
    #where did this come from?
```

To Optimize Betas (all weights within LSTM cells):

Stochastic Gradient Descent (SGD)

-- optimize over one sample each iteration

Mini-Batch SDG:

--optimize over b samples each iteration

Final Cost Function: $J^{(t)} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{|\mathcal{V}|} y_{i,j}^{(t)} \log \hat{y}_{i,j}^{(t)}$ -- "cross entropy error"

RNN-Based Language Models

Take-Aways

- Simple RNNs are powerful models but they are difficult to train:
 - Just two functions $h_{(t)}$ and $y_{(t)}$ where $h_{(t)}$ is a combination of $h_{(t-1)}$ and $x_{(t)}$.
 - Exploding and vanishing gradients make training difficult to converge.
- LSTM and GRU cells solve
 - Hidden states pass from one time-step to the next, allow for long-distance dependencies.
 - Gates are used to keep hidden states from changing rapidly (and thus keeps gradients under control).
 - To train: mini-batch stochastic gradient descent over cross-entropy cost

