

***Problem statement: Use an RNN to
predict the sentiment of the following
document***

I love cats → 1 (positive)

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1

Compute Hidden State
Contribution From Input

0	2	3
---	---	---

Input Token "I"
1 x 3
(1 token x E)

X

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

=

-2	-1
----	----

Proposed Hidden
State From Input
1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

0	2	3
---	---	---

Input Token "I"
1 x 3
(1 token x E)

X

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

=

-2	-1
----	----

Proposed Hidden State From Input
1 x 2

2 Compute Hidden State Contribution From Previous Hidden State

1	2
---	---

Previous Hidden State
1 x 2

X

3	1
-1	-2

HH Weights
2 x 2
(H x H)

=

1	2
---	---

Proposed Hidden State From Prev Hidden State
1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

0

2

3

Input Token "I"

1 x 3

(1 token x E)

X

3

1

-1

-2

0

1

EH Weights

3 x 2

(E x H)

=

-2

-1

Proposed Hidden State From Input

1 x 2

2 Compute Hidden State Contribution From Previous Hidden State

1

2

Previous Hidden State

1 x 2

X

3

1

-1

-2

HH Weights

2 x 2

(H x H)

=

1

2

Proposed Hidden State From Prev Hidden State

1 x 2

3 Update Hidden State

-2

-1

Proposed Hidden State From Input

+

1

2

Proposed Hidden State From Prev Hidden State

) =

-0.76

0.76

~

-1

1

New Hidden State

1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

0	2	3
---	---	---

Input Token "I"
1 x 3
(1 token x E)

X

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

=

-2	-1
----	----

Proposed Hidden State From Input
1 x 2

2 Compute Hidden State Contribution From Previous Hidden State

1	1
---	---

Previous Hidden State
1 x 2

X

3	1
-1	-2

HH Weights
2 x 2
(H x H)

=

1	2
---	---

Proposed Hidden State From Prev Hidden State
1 x 2

3 Update Hidden State

tanh(

-2	-1
----	----

Proposed Hidden State From Input

+

1	2
---	---

Proposed Hidden State From Prev Hidden State

) =

-0.76	0.76
-------	------

~

-1	1
----	---

New Hidden State
1 x 2

4 Compute Output

-1	1
----	---

New Hidden State
1 x 2

X

2	-2
---	----

HY Weights
2 x 1
(H x Y)

=

-4

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

Input Token "I" 1×3
(1 token \times E)

0	2	3
---	---	---

\times

3	1
-1	-2
0	1

EH Weights 3×2
(E \times H)

$=$

-2	-1
----	----

Proposed Hidden State From Input 1×2

4 Compute Output

-1	1
----	---

New Hidden State 1×2

\times

2	-2
---	----

HY Weights 2×1
(H \times Y)

$=$

-4

Sequence Step 1/4

2 Compute Hidden State Contribution From Previous Hidden State

1	2
---	---

Previous Hidden State 1×2

\times

3	1
-1	-2

HH Weights 2×2
(H \times H)

$=$

1	2
---	---

Proposed Hidden State From Prev Hidden State 1×2

3 Update Hidden State

$\tanh($

-2	-1
----	----

Proposed Hidden State From Input

$+$

1	2
---	---

Proposed Hidden State From Prev Hidden State

$) =$

-0.76	0.76
-------	------

\sim

-1	1
----	---

New Hidden State 1×2

5 Compute Output

$\text{sigmoid}($

-4

$) =$

0.18

Our predicted sentiment at sequence step 1

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1

Compute Hidden State
Contribution From Input

2	-1	2
---	----	---

Input Token
"Love"
1 x 3
(1 token x E)

X

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

=

7	6
---	---

Proposed Hidden
State From Input
1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

2	-1	2
---	----	---

Input Token
"Love"
1 x 3
(1 token x E)

X

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

=

7	6
---	---

Proposed Hidden State From Input
1 x 2

2 Compute Hidden State Contribution From Previous Hidden State

-1	1
----	---

Previous Hidden State
1 x 2

X

3	1
-1	-2

HH Weights
2 x 2
(H x H)

=

-4	-3
----	----

Proposed Hidden State From Prev Hidden State
1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

2	-1	2
---	----	---

Input Token
"Love"
1 x 3
(1 token x E)

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

X

7	6
---	---

Proposed Hidden State From Input
1 x 2

=

2 Compute Hidden State Contribution From Previous Hidden State

-1	1
----	---

Previous Hidden State
1 x 2

X

3	1
-1	-2

HH Weights
2 x 2
(H x H)

=

-4	-3
----	----

Proposed Hidden State From Prev Hidden State
1 x 2

3 Update Hidden State

tanh(

7	6
---	---

Proposed Hidden State From Input

+

-4	-3
----	----

Proposed Hidden State From Prev Hidden State

) =

.99	.99
-----	-----

~

1	1
---	---

New Hidden State
1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

Input Token "Love" 1×3 (1 token \times E)

2	-1	2
---	----	---

\times

3	1
-1	-2
0	1

EH Weights 3×2 (E \times H)

$=$

7	6
---	---

Proposed Hidden State From Input 1×2

2 Compute Hidden State Contribution From Previous Hidden State

-1	1
----	---

Previous Hidden State 1×2

\times

3	1
-1	-2

HH Weights 2×2 (H \times H)

$=$

-4	-3
----	----

Proposed Hidden State From Prev Hidden State 1×2

3 Update Hidden State

$\tanh($

7	6
---	---

Proposed Hidden State From Input

$+$

-4	-3
----	----

Proposed Hidden State From Prev Hidden State

$) =$

.99	.99
-----	-----

\sim

1	1
---	---

New Hidden State 1×2

4 Compute Output

1	1
---	---

New Hidden State 1×2

\times

2	-2
---	----

HY Weights 2×1 (H \times Y)

$=$

0

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

Input Token "Love" 1×3 (1 token \times E)

2	-1	2
---	----	---

 \times

3	1
-1	-2
0	1

EH Weights 3×2 (E \times H)

 $=$

7	6
---	---

Proposed Hidden State From Input 1×2

4 Compute Output

1	1
---	---

New Hidden State 1×2

 \times

2	-2
---	----

HY Weights 2×1 (H \times Y)

 $=$

0

Sequence Step 2/4

2 Compute Hidden State Contribution From Previous Hidden State

-1	1
----	---

Previous Hidden State 1×2

 \times

3	1
-1	-2

HH Weights 2×2 (H \times H)

 $=$

-4	-3
----	----

Proposed Hidden State From Prev Hidden State 1×2

3 Update Hidden State

$\tanh($

7	6
---	---

Proposed Hidden State From Input

 $+$

-4	-3
----	----

Proposed Hidden State From Prev Hidden State

 $=$

.99	.99
-----	-----

\sim

1	1
---	---

New Hidden State 1×2

5 Compute Output

$\text{sigmoid}($

0

 $=$

0.5

Our predicted sentiment at sequence step 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1

Compute Hidden State Contribution From Input

-1	-1	-2
----	----	----

Input Token

"Cats"

1 x 3

(1 token x E)

X

3	1
-1	-2
0	1

EH Weights

3 x 2

(E x H)

=

-2	3
----	---

Proposed Hidden State From Input

1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

-1	-1	-2
----	----	----

Input Token
"Cats"
1 x 3
(1 token x E)

X

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

=

-2	3
----	---

Proposed Hidden State From Input
1 x 2

2 Compute Hidden State Contribution From Previous Hidden State

1	1
---	---

Previous Hidden State
1 x 2

X

3	1
-1	-2

HH Weights
2 x 2
(H x H)

=

2	-1
---	----

Proposed Hidden State From Prev Hidden State
1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

-1

-1

-2

3

1

-1

-2

0

1

-2

3

Input Token
"Cats"
1 x 3
(1 token x E)

EH Weights
3 x 2
(E x H)

Proposed Hidden State From Input
1 x 2

2 Compute Hidden State Contribution From Previous Hidden State

1

1

3

1

-1

-2

2

-1

Previous Hidden State
1 x 2

HH Weights
2 x 2
(H x H)

Proposed Hidden State From Prev Hidden State
1 x 2

3 Update Hidden State

-2

3

2

-1

0

.96

Proposed Hidden State From Input

Proposed Hidden State From Prev Hidden State

0

1

New Hidden State
1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

-1	-1	-2
----	----	----

 \times

3	1
-1	-2
0	1

 =

-2	3
----	---

Input Token "Cats" 1×3 (1 token \times E)
EH Weights 3×2 (E \times H)
Proposed Hidden State From Input 1×2

2 Compute Hidden State Contribution From Previous Hidden State

1	1
---	---

 \times

3	1
-1	-2

 =

2	-1
---	----

Previous Hidden State 1×2
HH Weights 2×2 (H \times H)
Proposed Hidden State From Prev Hidden State 1×2

3 Update Hidden State

$\tanh($

-2	3
----	---

 $+$

2	-1
---	----

 $) =$

0	.96
---	-----

Proposed Hidden State From Input
Proposed Hidden State From Prev Hidden State
 \sim
New Hidden State 1×2

4 Compute Output

0	1
---	---

 \times

2	-2
---	----

 =

-2

New Hidden State 1×2
HY Weights 2×1 (H \times Y)

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

-1	-1	-2
----	----	----

 \times

3	1
-1	-2
0	1

 =

-2	3
----	---

Input Token "Cats" 1×3 (1 token \times E)
EH Weights 3×2 (E \times H)
Proposed Hidden State From Input 1×2

4 Compute Output

0	1
---	---

 \times

2	-2
---	----

 =

-2

New Hidden State 1×2
HY Weights 2×1 (H \times Y)

Sequence Step 3/4

2 Compute Hidden State Contribution From Previous Hidden State

1	1
---	---

 \times

3	1
-1	-2

 =

2	-1
---	----

Previous Hidden State 1×2
HH Weights 2×2 (H \times H)
Proposed Hidden State From Prev Hidden State 1×2

3 Update Hidden State

$\tanh($

-2	3
----	---

 $+$

2	-1
---	----

 $) =$

0	.96
---	-----

Proposed Hidden State From Input
Proposed Hidden State From Prev Hidden State

\sim

0	1
---	---

New Hidden State 1×2

5 Compute Output

$\text{sigmoid}($

-2

 $) =$

0.12

Our predicted sentiment at sequence step 3

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1

Compute Hidden State Contribution From Input

0	0	0
---	---	---

Input Token
"EMPTY_TOKEN"
1 x 3
(1 token x E)

X

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

=

0	0
---	---

Proposed Hidden State From Input
1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

0	0	0
---	---	---

Input Token
"EMPTY_TOKEN"
1 x 3
(1 token x E)

X

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

=

0	0
---	---

Proposed Hidden State From Input
1 x 2

2 Compute Hidden State Contribution From Previous Hidden State

0	1
---	---

Previous Hidden State
1 x 2

X

3	1
-1	-2

HH Weights
2 x 2
(H x H)

=

-1	-2
----	----

Proposed Hidden State From Prev Hidden State
1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

0	0	0
---	---	---

Input Token
"EMPTY_TOKEN"
1 x 3
(1 token x E)

X

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

=

0	0
---	---

Proposed Hidden State From Input
1 x 2

2 Compute Hidden State Contribution From Previous Hidden State

0	1
---	---

Previous Hidden State
1 x 2

X

3	1
-1	-2

HH Weights
2 x 2
(H x H)

=

-1	-2
----	----

Proposed Hidden State From Prev Hidden State
1 x 2

3 Update Hidden State

tanh(

0	0
---	---

Proposed Hidden State From Input

+

-1	-2
----	----

Proposed Hidden State From Prev Hidden State

) =

-.76	-.96
------	------

~

-1	-1
----	----

New Hidden State
1 x 2

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

0	0	0
---	---	---

Input Token
"EMPTY_TOKEN"
1 x 3
(1 token x E)

X

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

=

0	0
---	---

Proposed Hidden State From Input
1 x 2

2 Compute Hidden State Contribution From Previous Hidden State

0	1
---	---

Previous Hidden State
1 x 2

X

3	1
-1	-2

HH Weights
2 x 2
(H x H)

=

-1	-2
----	----

Proposed Hidden State From Prev Hidden State
1 x 2

3 Update Hidden State

0	0
---	---

Proposed Hidden State From Input

+

-1	-2
----	----

Proposed Hidden State From Prev Hidden State

) =

-0.76	-0.96
-------	-------

~

-1	-1
----	----

New Hidden State
1 x 2

4 Compute Output

-1	-1
----	----

New Hidden State
1 x 2

X

2	-2
---	----

HY Weights
2 x 1
(H x Y)

=

0

(S) Sequence Length = 4
(E) Embedding Size = 3
(H) Hidden State Dimensions = 2
(Y) Output Dimension = 1

1 Compute Hidden State Contribution From Input

Input Token "EMPTY_TOKEN" 1×3 (1 token \times E)

0	0	0
---	---	---

\times

3	1
-1	-2
0	1

EH Weights 3×2 (E \times H)

$=$

0	0
---	---

Proposed Hidden State From Input 1×2

4 Compute Output

-1	-1
----	----

New Hidden State 1×2

\times

2	-2
---	----

HY Weights 2×1 (H \times Y)

$=$

0

Sequence Step 4/4

2 Compute Hidden State Contribution From Previous Hidden State

0	1
---	---

Previous Hidden State 1×2

\times

3	1
-1	-2

HH Weights 2×2 (H \times H)

$=$

-1	-2
----	----

Proposed Hidden State From Prev Hidden State 1×2

3 Update Hidden State

$\tanh($

0	0
---	---

Proposed Hidden State From Input

$+$

-1	-2
----	----

Proposed Hidden State From Prev Hidden State

$) =$

-.76	-.96
------	------

\sim

-1	-1
----	----

New Hidden State 1×2

5 Compute Output

$\text{sigmoid}($

0

$) =$

0.5

Our predicted sentiment at sequence step 4

predicted y = 0.5
true y = 1

Binary Cross-Entropy (Log Loss) Function

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

Loss is $-\log(0.5) \rightarrow 0.3$

predicted y = 0.5
true y = 1

Binary Cross-Entropy (Log Loss) Function

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

Loss is $-\log(0.5) \rightarrow 0.3$

If this was the only sample in our training dataset, we would be done with 1 epoch. We then use our log loss function to calculate the partial derivatives needed to update our weights

3	1
-1	-2
0	1

EH Weights
3 x 2
(E x H)

3	1
-1	-2

HH Weights
2 x 2
(H x H)

2	-2
---	----

HY Weights
2 x 1
(H x Y)

predicted y = 0.78
true y = 1

Binary Cross-Entropy (Log Loss) Function

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

Loss is $-\log(0.78) \rightarrow 0.108$

Pretend we update our weights via backpropagation, and make a new prediction for “I love cats” of 0.78.

2	2
1	-1
-2	2

EH Weights
3 x 2
(**E** x **H**)

-1	2
-0	2

HH Weights
2 x 2
(**H** x **H**)

1	-3
---	----

HY Weights
2 x 1
(**H** x **Y**)

predicted y = 0.98
true y = 1

Binary Cross-Entropy (Log Loss) Function

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

Loss is $-\log(0.98) \rightarrow 0.008$

Pretend we update our weights via backpropagation, and make a new prediction for “I love cats” of 0.98.

3	2
2	-1
-2	2

EH Weights
3 x 2
(**E** x **H**)

-2	2
-0	2

HH Weights
2 x 2
(**H** x **H**)

1	-4
---	----

HY Weights
2 x 1
(**H** x **Y**)

Backpropagation Through Time

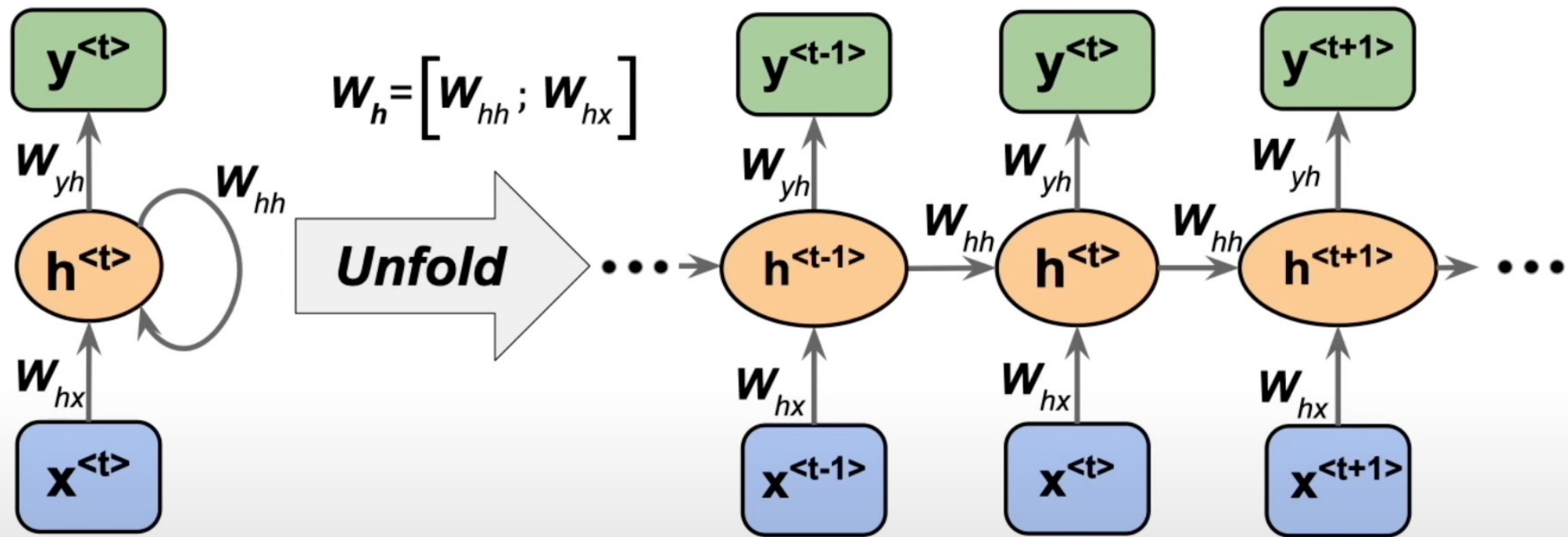
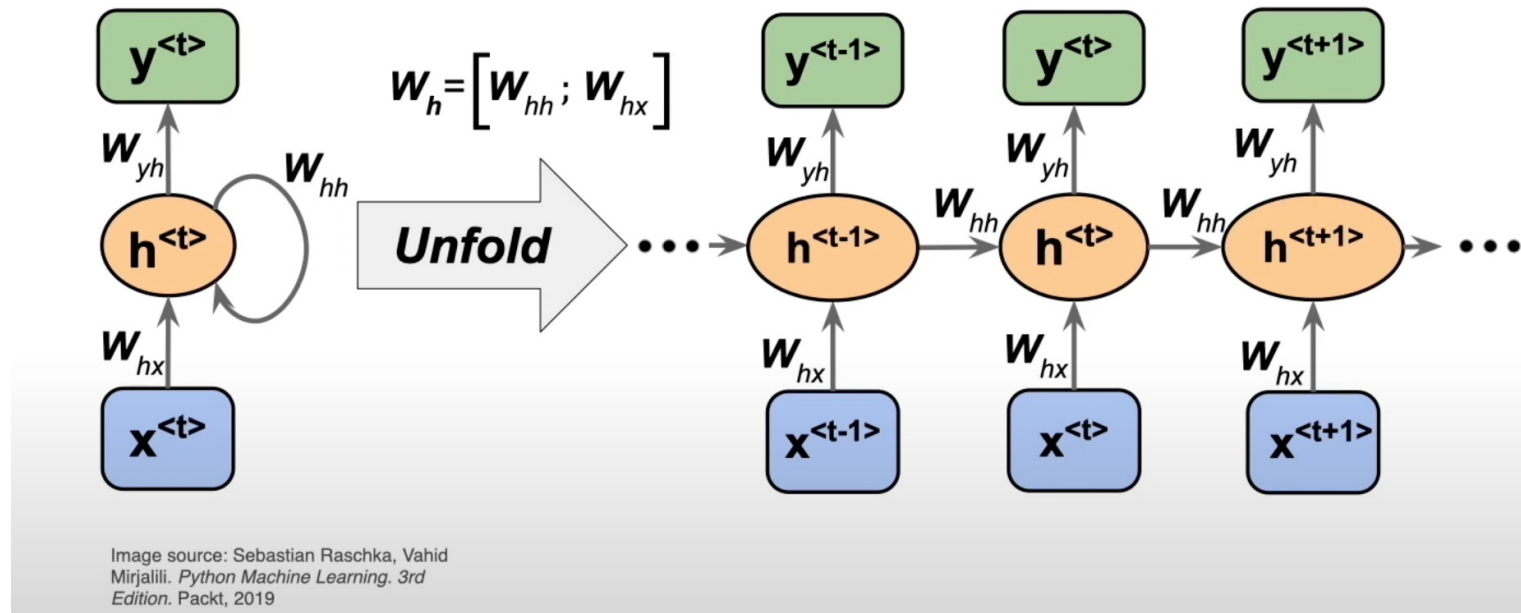


Image source: Sebastian Raschka, Vahid Mirjalili. *Python Machine Learning*. 3rd Edition. Packt, 2019

Backpropagation Through Time



$$L = \sum_{t=1}^T L^{(t)} \quad \frac{\partial L^{(t)}}{\partial \mathbf{W}_{hh}} = \frac{\partial L^{(t)}}{\partial y^{(t)}} \cdot \frac{\partial y^{(t)}}{\partial \mathbf{h}^{(t)}} \cdot \left(\sum_{k=1}^t \boxed{\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}}} \cdot \frac{\partial \mathbf{h}^{(k)}}{\partial \mathbf{W}_{hh}} \right)$$

computed as a multiplication of adjacent time steps:

This is very problematic:

Vanishing/Exploding gradient problem!

$$\frac{\partial \mathbf{h}^{(t)}}{\partial \mathbf{h}^{(k)}} = \prod_{i=k+1}^t \frac{\partial \mathbf{h}^{(i)}}{\partial \mathbf{h}^{(i-1)}}$$

Backpropagation Through Time

$$h_t = \sigma(wh_{t-1}).$$

Hidden state at sequence step t

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Definition of the sigmoid activation function

$$\sigma'(x) = \frac{d}{dx}\sigma(x) = \sigma(x)(1 - \sigma(x))$$

Derivative of the sigmoid activation function

$$\frac{\partial h_{t'}}{\partial h_t} = \prod_{k=1}^{t'-t} w \sigma'(wh_{t'-k})$$

Derivative of the hidden state

$$= \underbrace{w^{t'-t}}_{!!!} \prod_{k=1}^{t'-t} \sigma'(wh_{t'-k})$$