# Shortcomings of RNN

# Backward Propagation

- $\partial L/\partial W = (\partial L/\partial \hat{y}) \cdot (\partial \hat{y}/\partial H_3) \cdot (\partial g(z_3)/\partial z_3)[\, H_2 + W((\partial g(z_2)/\partial z_2)[\, H_1 + W((\partial g(z_1)/\partial z_1)[\, H_0 + W(\partial H_0/\partial W)])]\,)]$

- $\partial L/\partial U = (\partial L/\partial \hat{y}) \cdot (\partial \hat{y}/\partial H_3) \cdot (\partial g(z_3)/\partial z_3)[\, x_3 + (W\,(\partial g(z_2)/\partial z_2) \cdot [\, x_2 + (W\,(\partial g(z_1)/\partial z_1) \cdot [\, x_1 + (\partial WH_0/\partial U)\,])\,])\,]$

# Shortcomings of RNN

- **∂L/∂W**

- **∂L/∂U**

L = Loss

$\hat{y}$

$$H_0 \cdots\cdots\rightarrow H_3 \xrightarrow{W} H_2 \quad \cdots \quad \xrightarrow{W} H_n \xrightarrow{V} \hat{y}$$

$$x_1 \xrightarrow{U} H_3 \qquad x_2 \xrightarrow{U} H_2 \qquad x_3 \xrightarrow{U} H_n$$

# Shortcomings of RNN

- **∂L/∂W** → $\partial H_n / \partial W$

- **∂L/∂U** → $\partial H_n / \partial U$

L = Loss

$\hat{y}$

# Shortcomings of RNN

- **∂L/∂W** $\rightarrow \partial H_n/\partial W \rightarrow \partial H_{n-1}/\partial W$

- **∂L/∂U** $\rightarrow \partial H_n/\partial U \rightarrow \partial H_{n-1}/\partial U$

L = Loss

$\hat{y}$

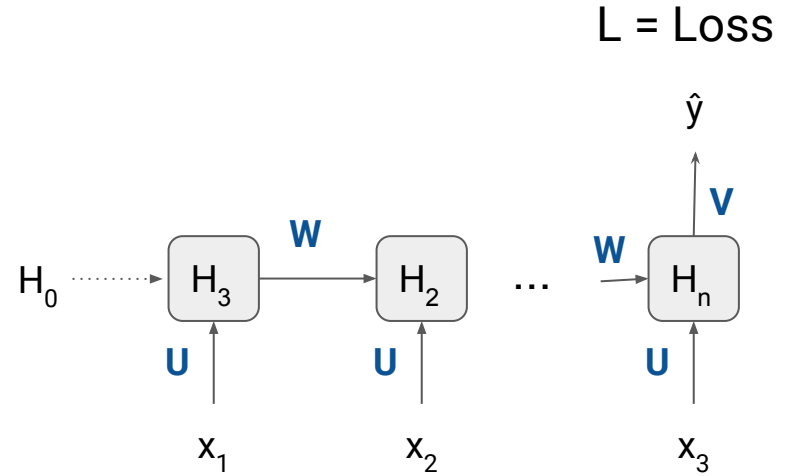# Shortcomings of RNN

- $\partial L/\partial W \rightarrow \partial H_n/\partial W \rightarrow \partial H_{n-1}/\partial W \ldots \partial H_0/\partial W$

- $\partial L/\partial U \rightarrow \partial H_n/\partial U \rightarrow \partial H_{n-1}/\partial U \ldots \partial H_0/\partial U$

L = Loss

$\hat{y}$

$H_0 \cdots\!\rightarrow \boxed{H_3} \xrightarrow{\ W\ } \boxed{H_2} \quad \ldots \quad \xrightarrow{\ W\ } \boxed{H_n} \rightarrow$

V

U          U          U

$x_1$          $x_2$          $x_3$

# Vanishing Gradient

- **∂L/∂W** → $\partial H_n/\partial W$ → $\partial H_{n-1}/\partial W$ ... $\partial H_0/\partial W$

- **∂L/∂U** → $\partial H_n/\partial U$ → $\partial H_{n-1}/\partial U$ ... $\partial H_0/\partial U$

L = Loss

$\hat{y}$

**V**

$H_0$ ·······→ $H_3$ **W** → $H_2$ ... **W** → $H_n$

**U** **U** **U**

$x_1$ $x_2$ $x_3$

Analytics Vidhya
Learn everything about analytics

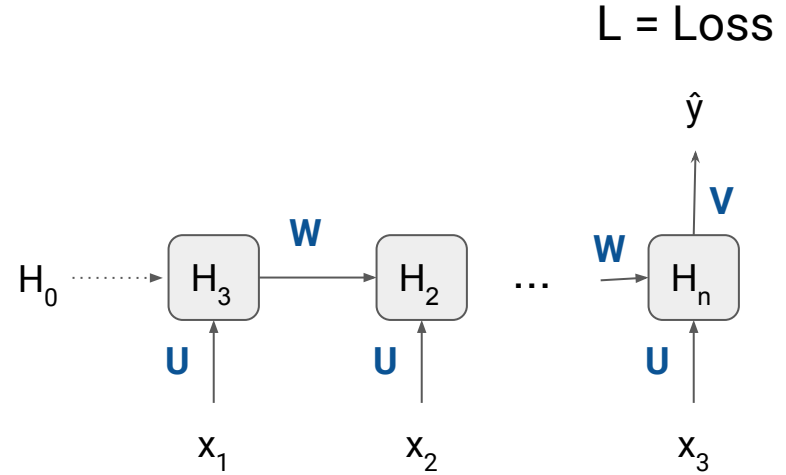# Vanishing Gradient

- $\partial L/\partial W \rightarrow \partial H_n/\partial W \rightarrow \partial H_{n-1}/\partial W \ldots \partial H_0/\partial W$

- $\partial L/\partial U \rightarrow \partial H_n/\partial U \rightarrow \partial H_{n-1}/\partial U \ldots \partial H_0/\partial U$

- If gradients < 1

L = Loss

$\hat{y}$

V

$H_0$ ┈┈→ $H_3$ —W→ $H_2$ … —W→ $H_n$
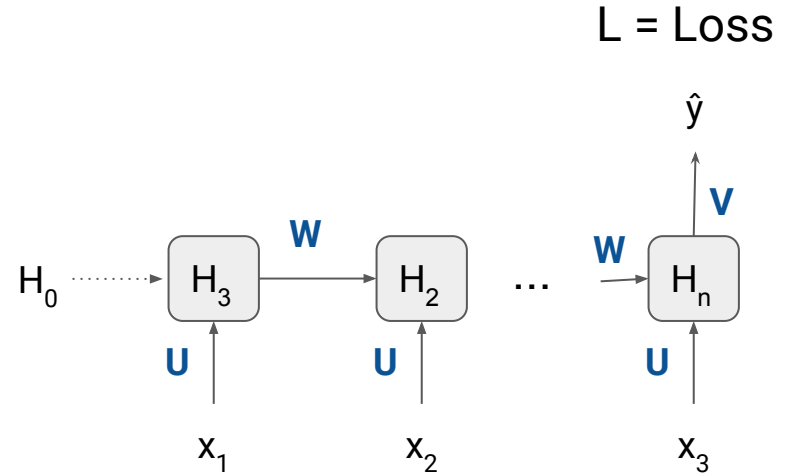
U          U          U

$x_1$        $x_2$        $x_3$

# Vanishing Gradient
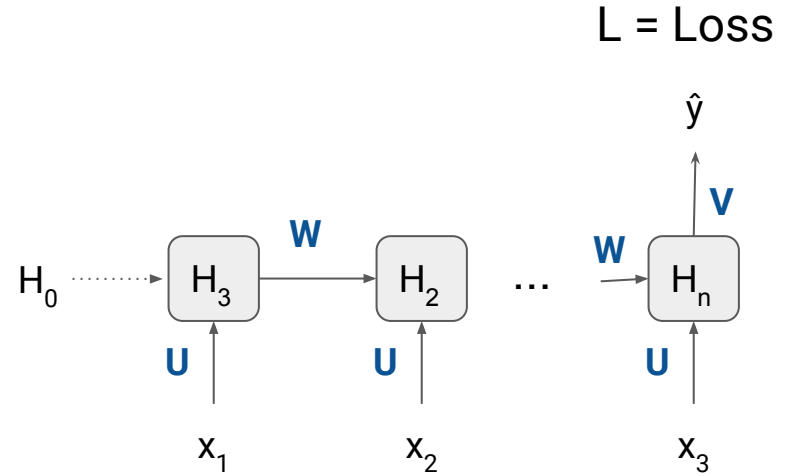
- $\partial L/\partial W \rightarrow \partial H_n/\partial W \rightarrow \partial H_{n-1}/\partial W \ldots \partial H_0/\partial W$

- $\partial L/\partial U \rightarrow \partial H_n/\partial U \rightarrow \partial H_{n-1}/\partial U \ldots \partial H_0/\partial U$

- If gradients < 1

  then $\partial L/\partial W$ and $\partial L/\partial U$ will be infinitesimally small

L = Loss

$\hat{y}$

$H_0 \cdots\cdots\rightarrow \boxed{H_3} \xrightarrow{\ W\ } \boxed{H_2} \quad \cdots \quad \xrightarrow{\ W\ } \boxed{H_n}$

$U \qquad\qquad U \qquad\qquad V \quad U$

$x_1 \qquad\qquad x_2 \qquad\qquad x_3$

Analytics Vidhya
Learn everything about analytics

# Vanishing Gradient

- $\partial L/\partial W \rightarrow \partial H_n/\partial W \rightarrow \partial H_{n-1}/\partial W \dots \partial H_0/\partial W$

- $\partial L/\partial U \rightarrow \partial H_n/\partial U \rightarrow \partial H_{n-1}/\partial U \dots \partial H_0/\partial U$

- If gradients < 1

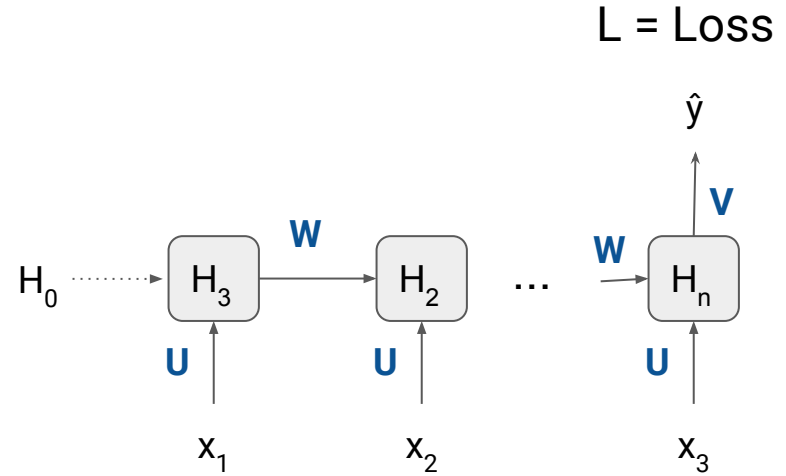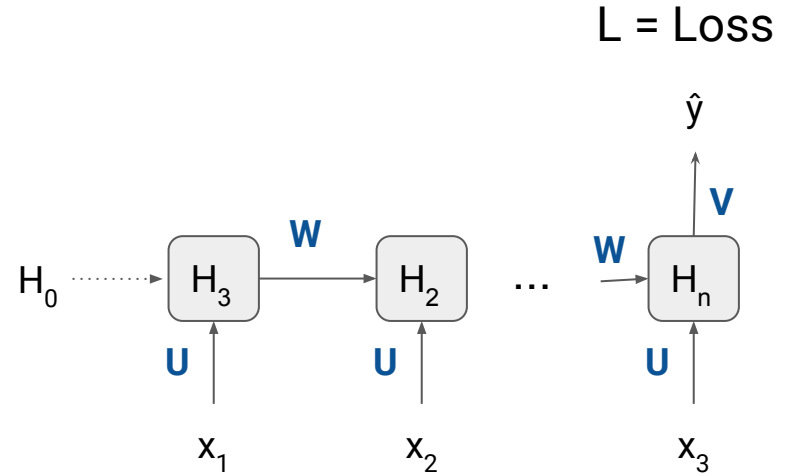  then $\partial L/\partial W$ and $\partial L/\partial U$ will be infinitesimally small

  $(0.5)^{10} = 0.00097$

L = Loss

$\hat{y}$

$H_0 \cdots \rightarrow [H_3] \xrightarrow{W} [H_2] \dots \xrightarrow{W} [H_n]$

$U \quad U \quad U$

$x_1 \quad x_2 \quad x_3$

V

# Vanishing Gradient

- $\partial L / \partial W \approx \partial L / \partial U \approx 0$

L = Loss

$\hat{y}$

**V**

$H_0$ ┈┈> [ $H_3$ ] **W** > [ $H_2$ ] … **W** > [ $H_n$ ]

**U**      **U**      **U**

$x_1$      $x_2$      $x_3$

Analytics Vidhya
Learn everything about analytics

# Vanishing Gradient

- $\partial L/\partial W \approx \partial L/\partial U \approx 0$

- $W = W - \alpha \, (\partial L/\partial W)$
  $U = U - \alpha \, (\partial L/\partial U)$

**updating weights**

L = Loss

$\hat{y}$

$H_0 \cdots\cdots\rightarrow \boxed{H_3} \xrightarrow{\;W\;} \boxed{H_2} \;\cdots\; \xrightarrow{W} \boxed{H_n} \xrightarrow{V}$

U          U          U

$x_1$          $x_2$          $x_3$

Analytics Vidhya
Learn everything about analytics

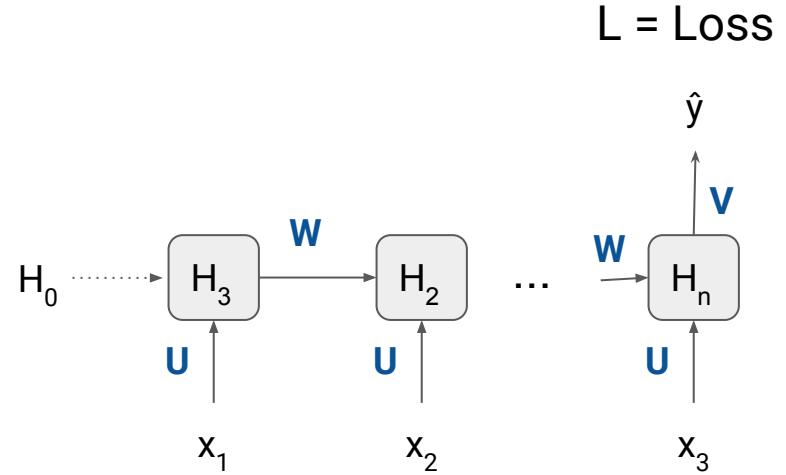# Vanishing Gradient
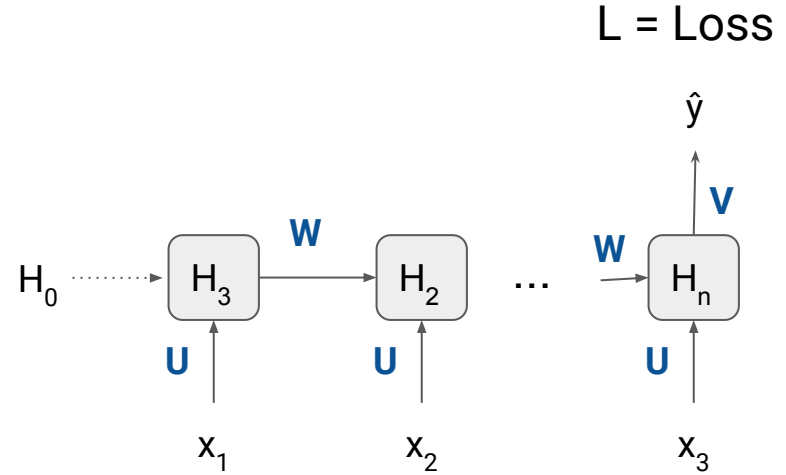
- $\partial L/\partial W \approx \partial L/\partial U \approx 0$

- $W = W - \alpha\,(\partial L/\partial W)$
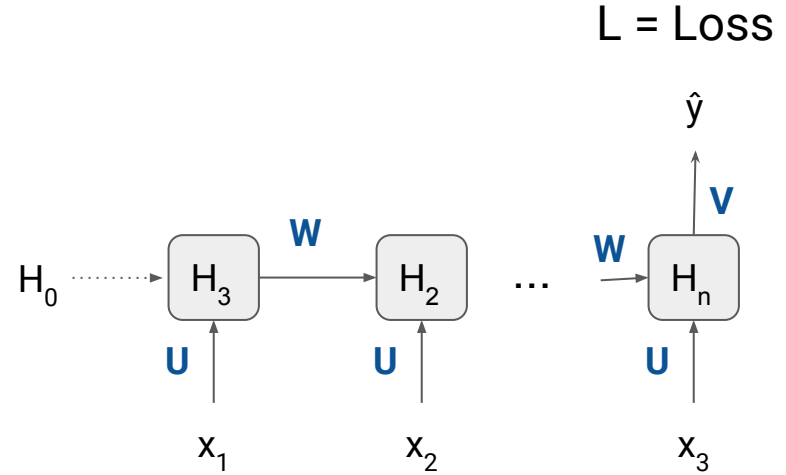  $U = U - \alpha\,(\partial L/\partial U)$

- $W_{before} \approx W_{after}$
  $U_{before} \approx U_{after}$

**updating weights**

L = Loss

# Exploding Gradient

- $\partial L/\partial W \rightarrow \partial H_n/\partial W \rightarrow \partial H_{n-1}/\partial W \dots \partial H_0/\partial W$

- $\partial L/\partial U \rightarrow \partial H_n/\partial U \rightarrow \partial H_{n-1}/\partial U \dots \partial H_0/\partial U$

L = Loss

$\hat{y}$

V

$H_0 \cdots\rightarrow \boxed{H_3} \xrightarrow{W} \boxed{H_2} \quad \dots \quad \xrightarrow{W} \boxed{H_n}$

U          U                  U

$x_1$        $x_2$              $x_3$

# Exploding Gradient
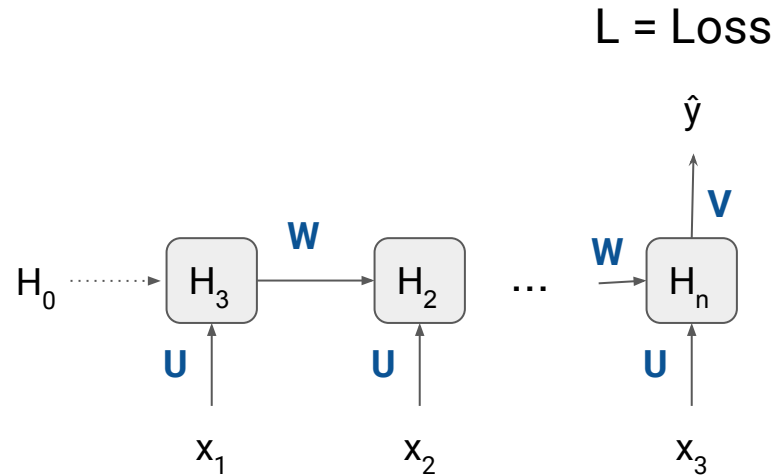
- $\partial L/\partial W \rightarrow \partial H_n/\partial W \rightarrow \partial H_{n-1}/\partial W \ldots \partial H_0/\partial W$

- $\partial L/\partial U \rightarrow \partial H_n/\partial U \rightarrow \partial H_{n-1}/\partial U \ldots \partial H_0/\partial U$

- If gradients > 1

  then $\partial L/\partial W$ and $\partial L/\partial U$ are very large

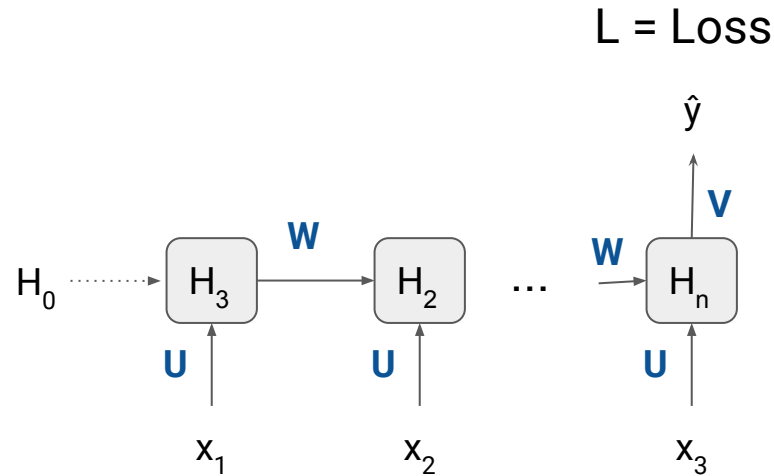L = Loss

$\hat{y}$

$V$

$H_0 \cdots\cdots\rightarrow \boxed{H_3} \xrightarrow{W} \boxed{H_2} \quad \ldots \quad \xrightarrow{W} \boxed{H_n}$

$U \quad\quad\quad U \quad\quad\quad U$

$x_1 \quad\quad\quad x_2 \quad\quad\quad x_3$

# Exploding Gradient

- **∂L/∂W** → $\partial H_n/\partial W$ → $\partial H_{n-1}/\partial W$ ... $\partial H_0/\partial W$

- **∂L/∂U** → $\partial H_n/\partial U$ → $\partial H_{n-1}/\partial U$ ... $\partial H_0/\partial U$

- If gradients > 1

  then ∂L/∂W and ∂L/∂U are very large

  $(1.5)^{10} = 57.665$

L = Loss

$\hat{y}$

$H_0$ ·····> $H_3$ --**W**--> $H_2$ ... **W** $H_n$
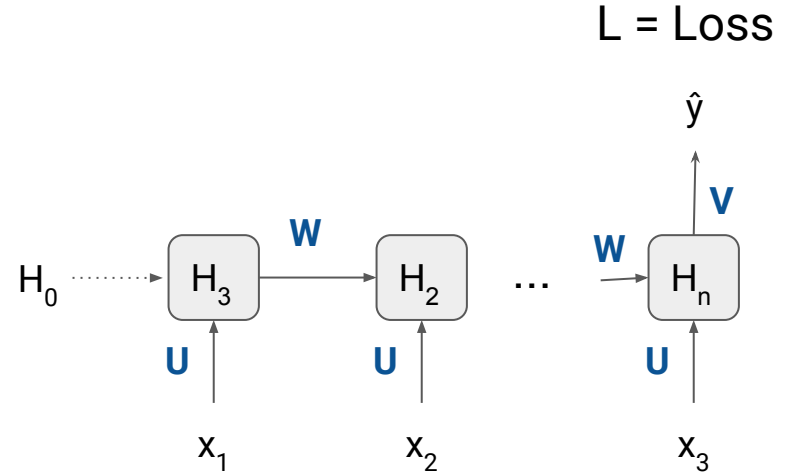
**V**

**U**   **U**   **U**

$x_1$   $x_2$   $x_3$

# Exploding Gradient

- $W = W - \alpha \, (\partial L / \partial W)$
  $U = U - \alpha \, (\partial L / \partial U)$

**updating weights**

L = Loss

$\hat{y}$

V

$H_0$ ·····> $H_3$ — **W** → $H_2$ ... **W** → $H_n$

U       U       U

$x_1$     $x_2$     $x_3$

# Issues with RNN

- **Vanishing Gradient**

    RNN does not perform well for long sentences

# Issues with RNN

- **Vanishing Gradient**

  RNN does not perform well for long sentences

  Eg: The writer of the books

# Issues with RNN

- **Vanishing Gradient**

   RNN does not perform well for long sentences

   Eg: The writer of the books

Analytics Vidhya
Learn everything about analytics

# Issues with RNN

- **Vanishing Gradient**

  RNN does not perform well for long sentences

  Eg: The writer of the books is

Analytics Vidhya
Learn everything about analytics

# Issues with RNN

- **Vanishing Gradient**

  RNN does not perform well for long sentences

  Eg: The writer of the books is
      The writer of the books are

# Issues with RNN

- **Vanishing Gradient**

  RNN does not perform well for long sentences

  Eg: The writer of the books is
       The writer of the books are

- **Exploding Gradient**

  - Gradients are large

# Issues with RNN

- **Vanishing Gradient**

  RNN does not perform well for long sentences

  Eg: The writer of the books is
       The writer of the books are

- **Exploding Gradient**

  - Gradients are large
  - Poor predictions

# Issues with RNN

How to mitigate exploding and vanishing gradients?

# Issues with RNN

How to mitigate exploding and vanishing gradients?

| Exploding Gradients |
|---|
| Gradient Clipping |

# Issues with RNN

How to mitigate exploding and vanishing gradients?

| Exploding Gradients |
|---|
| Gradient Clipping |

Threshold = 0.1

- if gradient > threshold, then gradient = threshold
- if gradient <= threshold, then gradient = gradient

Analytics Vidhya
Learn everything about analytics

# Issues with RNN

How to mitigate exploding and vanishing gradients?

| Exploding Gradients | Vanishing Gradients |
|---|---|
| Gradient Clipping | LSTM or GRU |

Analytics Vidhya
Learn everything about analytics

Thank You