# Recurrent Neural Networks

Amir H. Payberah
`payberah@kth.se`
07/12/2018
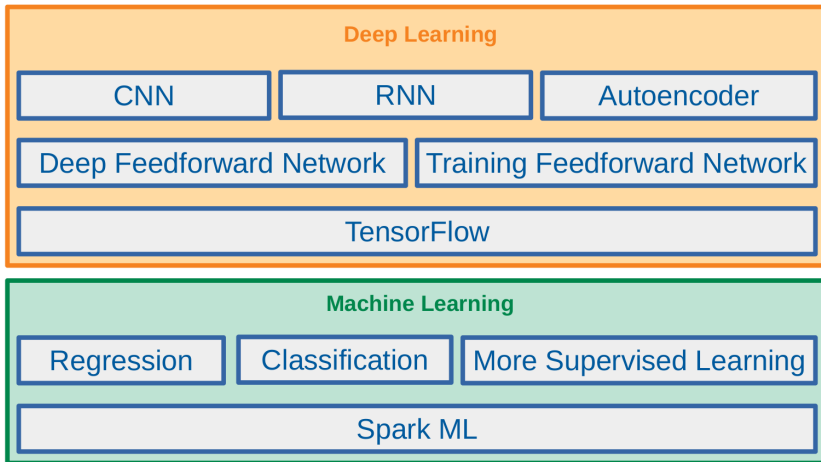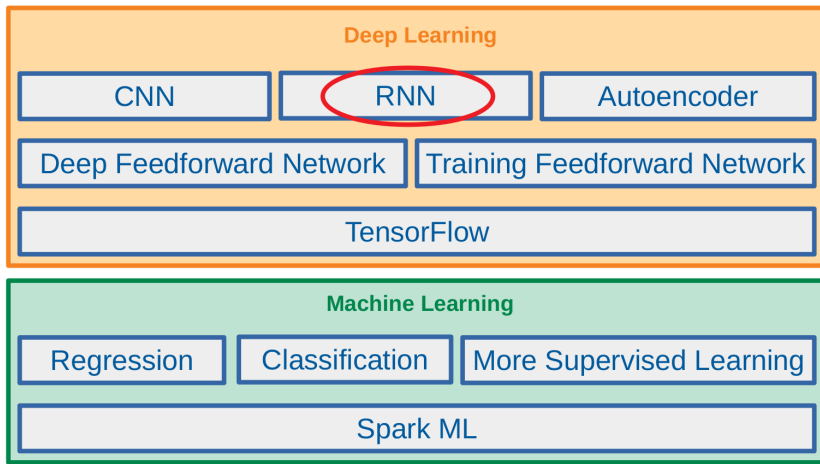
https://id2223kth.github.io

# Let's Start With An Example

- Language modeling is the task of predicting what word comes next.



the students opened their _____ → books, laptops, exams, minds

▶ More formally: given a sequence of words $x^{(1)}, x^{(2)}, \cdots, x^{(t)}$, compute the probability distribution of the next word $x^{(t+1)}$:

$$p(x^{(t+1)} = w_j | x^{(t)}, \cdots x^{(1)})$$

▶ More formally: given a sequence of words $x^{(1)}, x^{(2)}, \cdots, x^{(t)}$, compute the probability distribution of the next word $x^{(t+1)}$:

$$p(x^{(t+1)} = w_j | x^{(t)}, \cdots x^{(1)})$$

▶ $w_j$ is a word in vocabulary $V = \{w_1, \cdots, w_v\}$.

- `the students opened their ___`

- the students opened their ___

- How to learn a Language Model?

# n-gram Language Models

- `the students opened their ___`

- How to learn a Language Model?

- Learn a n-gram Language Model!

- the students opened their ___

- How to learn a Language Model?

- Learn a n-gram Language Model!

- A n-gram is a chunk of n consecutive words.

- the students opened their ___

- How to learn a Language Model?

- Learn a n-gram Language Model!

- A n-gram is a chunk of n consecutive words.
  - Unigrams: "the", "students", "opened", "their"

# n-gram Language Models

- `the students opened their ___`

- How to learn a Language Model?

- Learn a n-gram Language Model!

- A n-gram is a chunk of n consecutive words.
  - Unigrams: `"the"`, `"students"`, `"opened"`, `"their"`
  - Bigrams: `"the students"`, `"students opened"`, `"opened their"`

# n-gram Language Models

- `the students opened their ___`

- How to learn a Language Model?

- Learn a n-gram Language Model!

- A n-gram is a chunk of n consecutive words.
  - Unigrams: `"the"`, `"students"`, `"opened"`, `"their"`
  - Bigrams: `"the students"`, `"students opened"`, `"opened their"`
  - Trigrams: `"the students opened"`, `"students opened their"`

# n-gram Language Models

▶ the students opened their ___

▶ How to learn a Language Model?

▶ Learn a n-gram Language Model!

▶ A n-gram is a chunk of n consecutive words.
  • Unigrams: "the", "students", "opened", "their"
  • Bigrams: "the students", "students opened", "opened their"
  • Trigrams: "the students opened", "students opened their"
  • 4-grams: "the students opened their"

# n-gram Language Models

- `the students opened their ___`

- How to learn a Language Model?

- Learn a n-gram Language Model!

- A n-gram is a chunk of n consecutive words.
  - Unigrams: `"the"`, `"students"`, `"opened"`, `"their"`
  - Bigrams: `"the students"`, `"students opened"`, `"opened their"`
  - Trigrams: `"the students opened"`, `"students opened their"`
  - 4-grams: `"the students opened their"`

- Collect statistics about how frequent different n-grams are, and use these to predict next word.

- ▶ Suppose we are learning a 4-gram Language Model.
  - • $x^{(t+1)}$ depends only on the preceding 3 words $\{x^{(t)}, x^{(t-1)}, x^{(t-2)}\}$.

~~as the proctor started the clock, the~~ *students opened their* _____

discard                         condition on this

# n-gram Language Models - Example

- ▶ Suppose we are learning a 4-gram Language Model.
  - • $x^{(t+1)}$ depends only on the preceding 3 words $\{x^{(t)}, x^{(t-1)}, x^{(t-2)}\}$.

~~as the proctor started the clock, the~~ students opened their _____
        discard                          condition on this

$$p(w_j | \text{students opened their}) = \frac{\text{students opened their } w_j}{\text{students opened their}}$$

# n-gram Language Models - Example

- Suppose we are learning a 4-gram Language Model.
  - $x^{(t+1)}$ depends only on the preceding 3 words $\{x^{(t)}, x^{(t-1)}, x^{(t-2)}\}$.

~~as the proctor started the clock, the~~ *students opened their* ____
    discard                              condition on this

$$p(w_j | \text{students opened their}) = \frac{\text{students opened their } w_j}{\text{students opened their}}$$

- In the corpus:
  - `"students opened their"` occurred 1000 times

▶ Suppose we are learning a 4-gram Language Model.

- $x^{(t+1)}$ depends only on the preceding 3 words $\{x^{(t)}, x^{(t-1)}, x^{(t-2)}\}$.

~~as the proctor started the clock, the~~ students opened their _____
discard                                    condition on this

$$p(w_j|\text{students opened their}) = \frac{\text{students opened their } w_j}{\text{students opened their}}$$

▶ In the corpus:

- `"students opened their"` occurred 1000 times
- `"students opened their books` occurred 400 times:
  $p(\text{books}|\text{students opened their}) = 0.4$

# n-gram Language Models - Example

- Suppose we are learning a 4-gram Language Model.
  - $x^{(t+1)}$ depends only on the preceding 3 words $\{x^{(t)}, x^{(t-1)}, x^{(t-2)}\}$.

~~as the proctor started the clock, the~~ students opened their _____
discard          condition on this

$$p(w_j|\text{students opened their}) = \frac{\text{students opened their } w_j}{\text{students opened their}}$$

- In the corpus:
  - "students opened their" occurred 1000 times
  - "students opened their books occurred 400 times:
    $p(\text{books}|\text{students opened their}) = 0.4$
  - "students opened their exams occurred 100 times:
    $p(\text{exams}|\text{students opened their}) = 0.1$

$$p(w_j|\text{students opened their}) = \frac{\text{students opened their } w_j}{\text{students opened their}}$$

$$p(w_j | \text{students opened their}) = \frac{\text{students opened their } w_j}{\text{students opened their}}$$

▶ What if "students opened their $w_j$" never occurred in data? Then $w_j$ has probability 0!

$$p(\text{w}_j|\text{students opened their}) = \frac{\text{students opened their w}_j}{\text{students opened their}}$$

▶ What if `"students opened their w`$_j$`"` never occurred in data? Then `w`$_j$ has probability 0!

▶ What if `"students opened their"` never occurred in data? Then we can't calculate probability for any `w`$_j$!

$$p(w_j|\text{students opened their}) = \frac{\text{students opened their } w_j}{\text{students opened their}}$$

▶ What if "students opened their w_j" never occurred in data? Then w_j has probability 0!

▶ What if "students opened their" never occurred in data? Then we can't calculate probability for any w_j!

▶ Increasing n makes sparsity problems worse.
  • Typically we can't have n bigger than 5.

$$p(\mathtt{w_j}|\text{students opened their}) = \frac{\text{students opened their } \mathtt{w_j}}{\text{students opened their}}$$

$$p(\mathtt{w_j}|\text{students opened their}) = \frac{\text{students opened their } \mathtt{w_j}}{\text{students opened their}}$$

▶ For "`students opened their w_j`", we need to store count for all possible 4-grams.

▶ The model size is in the order of $\mathtt{O(exp(n))}$.

▶ Increasing $\mathtt{n}$ makes model size huge.

▶ Recall the Language Modeling task:
- Input: sequence of words $x^{(1)}, x^{(2)}, \cdots, x^{(t)}$
- Output: probability dist of the next word $p(x^{(t+1)} = w_j | x^{(t)}, \cdots, x^{(1)})$

▶ Recall the Language Modeling task:
  • Input: sequence of words $x^{(1)}, x^{(2)}, \cdots, x^{(t)}$
  • Output: probability dist of the next word $p(x^{(t+1)} = w_j | x^{(t)}, \cdots, x^{(1)})$

▶ One-Hot encoding
  • Represent a categorical variable as a binary vector.
  • All recodes are zero, except the index of the integer, which is one.
  • Each embedded word $\mathbf{e}^{(t)} = \mathbf{E}^\intercal x^{(t)}$ is a one-hot vector of size vocabulary size.

$$
\begin{array}{llll}
& & \text{opened} & \text{word V} \\
& \text{students} & & \\
\mathbf{x}^{(1)} & \text{students} = [1, & 0, & 0, \; 0, \; 0, \; 0, \; \ldots, \; 0] \\
\mathbf{x}^{(2)} & \text{opened} = [0, & 1, & 0, \; 0, \; 0, \; 0, \; \ldots, \; 0] \\
\mathbf{x}^{(3)} & \text{their} = [0, & 0, & 1, \; 0, \; 0, \; 0, \; \ldots, \; 0] \\
\mathbf{x}^{(4)} & \text{book} = [0, & 0, & 0, \; 1, \; 0, \; 0, \; \ldots, \; 0] \\
\end{array}
$$

$\mathbf{e}^{(t)}$

- A MLP model
  - **Input**: words $x^{(1)}, x^{(2)}, x^{(3)}, x^{(4)}$
  - **Input layer**: one-hot vectors $e^{(1)}, e^{(2)}, e^{(3)}, e^{(4)}$
  - **Hidden layer**: $h = f(w^{\intercal}e)$, $f$ is an activation function.
  - **Output**: $\hat{y} = \mathtt{softmax}(v^{\intercal}h)$

▶ Improvements over n-gram LM:
- No sparsity problem
- Model size is `O(n)` not `O(exp(n))`

- ▶ Improvements over n-gram LM:
  - • No sparsity problem
  - • Model size is `O(n)` not `O(exp(n))`

- ▶ Remaining problems:
  - • It is fixed 4 in our example, which is small
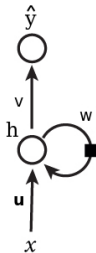  - • We need a neural architecture that can process any length input

# Recurrent Neural Networks (RNN)

- The idea behind Recurrent neural networks (RNN) is to make use of sequential data.

- The idea behind Recurrent neural networks (RNN) is to make use of sequential data.
  - Until here, we assume that all inputs (and outputs) are independent of each other.

- The idea behind Recurrent neural networks (RNN) is to make use of sequential data.
  - Until here, we assume that all inputs (and outputs) are independent of each other.
  - It is a bad idea for many tasks, e.g., predicting the next word in a sentence (it's better to know which words came before it).

- The idea behind Recurrent neural networks (RNN) is to make use of sequential data.
  - Until here, we assume that all inputs (and outputs) are independent of each other.
  - It is a bad idea for many tasks, e.g., predicting the next word in a sentence (it's better to know which words came before it).

- They can analyze time series data and predict the future.

- The idea behind Recurrent neural networks (RNN) is to make use of sequential data.
  - Until here, we assume that all inputs (and outputs) are independent of each other.
  - It is a bad idea for many tasks, e.g., predicting the next word in a sentence (it's better to know which words came before it).

- They can analyze time series data and predict the future.

- They can work on sequences of arbitrary lengths, rather than on fixed-sized inputs.

- Neurons in an RNN have connections pointing backward.
- RNNs have memory, which captures information about what has been calculated so far.

- Unfolding the network: represent a network against the time axis.
  - We write out the network for the complete sequence.

- Unfolding the network: represent a network against the time axis.
  - We write out the network for the complete sequence.
- For example, if the sequence we care about is a sentence of three words, the network would be unfolded into a 3-layer neural network.
  - One layer for each word.

▶ $h^{(t)} = f(\mathbf{u}^\mathsf{T} \mathbf{x}^{(t)} + wh^{(t-1)})$, where $f$ is an activation function, e.g., `tanh` or `ReLU`.

- $h^{(t)} = f(\mathbf{u}^{\mathsf{T}}\mathbf{x}^{(t)} + wh^{(t-1)})$, where `f` is an activation function, e.g., `tanh` or `ReLU`.
- $\hat{y}^{(t)} = g(vh^{(t)})$, where `g` can be the `softmax` function.

- $h^{(t)} = f(\mathbf{u}^{\mathsf{T}} \mathbf{x}^{(t)} + w h^{(t-1)})$, where $f$ is an activation function, e.g., `tanh` or `ReLU`.

- $\hat{y}^{(t)} = g(v h^{(t)})$, where $g$ can be the `softmax` function.

- $\text{cost}(y^{(t)}, \hat{y}^{(t)}) = \text{cross\_entropy}(y^{(t)}, \hat{y}^{(t)}) = -\sum y^{(t)} \log \hat{y}^{(t)}$

- $y^{(t)}$ is the correct word at time step $t$, and $\hat{y}^{(t)}$ is the prediction.

▶ Each recurrent neuron has three sets of weights: **u**, **w**, and **v**.

▶ **u**: the weights for the inputs $\mathbf{x}^{(t)}$.

- **u**: the weights for the inputs $\mathbf{x}^{(t)}$.

- $\mathbf{x}^{(t)}$: is the input at time step $t$.

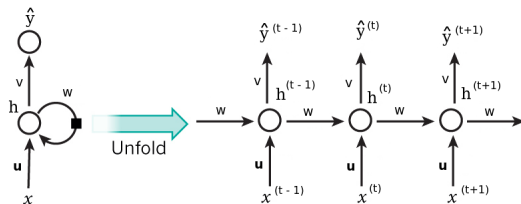- For example, $\mathbf{x}^{(1)}$ could be a one-hot vector corresponding to the first word of a sentence.

▶ w: the weights for the hidden state of the previous time step $h^{(t-1)}$.

- ▶ $w$: the weights for the hidden state of the previous time step $h^{(t-1)}$.
- ▶ $h^{(t)}$: is the hidden state (memory) at time step $t$.
  - $h^{(t)} = \tanh(\mathbf{u}^\mathsf{T}\mathbf{x}^{(t)} + wh^{(t-1)})$
  - $h^{(0)}$ is the initial hidden state.

▶ v: the weights for the hidden state of the current time step $h^{(t)}$.

- $v$: the weights for the hidden state of the current time step $h^{(t)}$.
- $\hat{\mathbf{y}}^{(t)}$ is the output at step $t$.
- $\hat{\mathbf{y}}^{(t)} = \texttt{softmax}(vh^{(t)})$

- ▶ $v$: the weights for the hidden state of the current time step $h^{(t)}$.
- ▶ $\hat{\mathbf{y}}^{(t)}$ is the output at step $t$.
- ▶ $\hat{\mathbf{y}}^{(t)} = \text{softmax}(vh^{(t)})$
- ▶ For example, if we wanted to predict the next word in a sentence, it would be a vector of probabilities across our vocabulary.
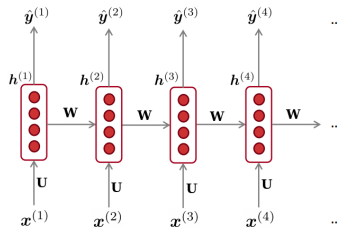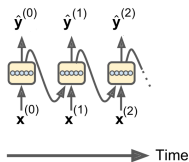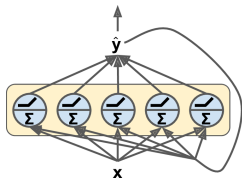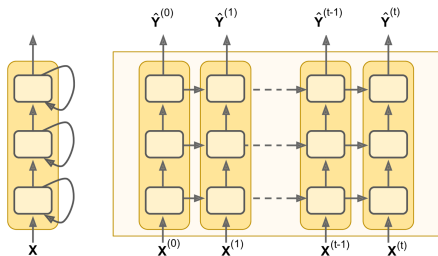
▶ At each time step $t$, every neuron of a layer receives both the input vector $\mathbf{x}^{(t)}$ and the output vector from the previous time step $\mathbf{h}^{(t-1)}$.

$$\mathbf{h}^{(t)} = \tanh(\mathbf{u}^\mathsf{T}\mathbf{x}^{(t)} + \mathbf{w}^\mathsf{T}\mathbf{h}^{(t-1)})$$
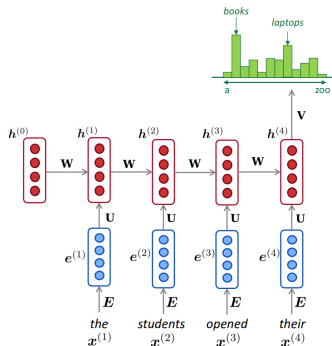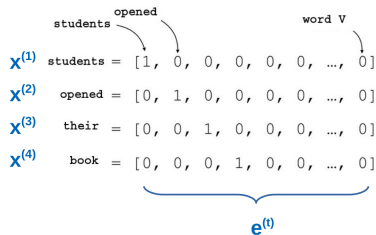$$\mathbf{y}^{(t)} = \texttt{sigmoid}(\mathbf{v}^\mathsf{T}\mathbf{h}^{(t)})$$

- Stacking multiple layers of cells gives you a deep RNN.

# Let's Back to Language Model Example
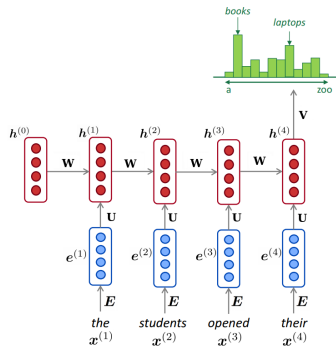
▶ The input **x** will be a sequence of words (each $\mathbf{x}^{(t)}$ is a single word).

▶ Each embedded word $\mathbf{e}^{(t)} = \mathbf{E}^{\mathsf{T}}\mathbf{x}^{(t)}$ is a one-hot vector of size vocabulary size.

▶ Let's recap the equations for the RNN:
- $\mathbf{h}^{(t)} = \tanh(\mathbf{u}^{\mathsf{T}}\mathbf{e}^{(t)} + \mathbf{w}\mathbf{h}^{(t-1)})$
- $\hat{\mathbf{y}}^{(t)} = \mathtt{softmax}(\mathbf{v}\mathbf{h}^{(t)})$

▶ Let's recap the equations for the RNN:
  • $h^{(t)} = \tanh(\mathbf{u}^\mathsf{T}\mathbf{e}^{(t)} + wh^{(t-1)})$
  • $\hat{\mathbf{y}}^{(t)} = \texttt{softmax}(vh^{(t)})$

▶ The output $\hat{\mathbf{y}}^{(t)}$ is a vector of vocabulary size elements.

- Let's recap the equations for the RNN:
  - $h^{(t)} = \tanh(\mathbf{u}^\intercal \mathbf{e}^{(t)} + wh^{(t-1)})$
  - $\hat{\mathbf{y}}^{(t)} = \texttt{softmax}(vh^{(t)})$

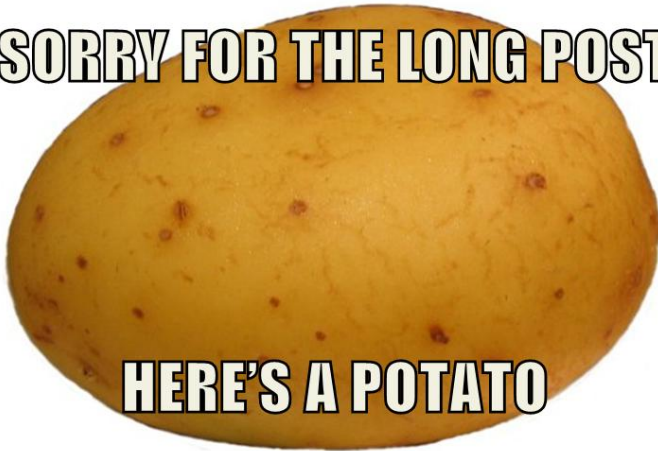- The output $\hat{\mathbf{y}}^{(t)}$ is a vector of vocabulary size elements.

- Each element of $\hat{\mathbf{y}}^{(t)}$ represents the probability of that word being the next word in the sentence.

SORRY FOR THE LONG POST

HERE'S A POTATO

# RNN in TensorFlow

▶ Manul implementation of an RNN

```python
# make the dataset
n_inputs = 3
n_neurons = 5

X0_batch = np.array([[0, 1, 2], [3, 4, 5], [6, 7, 8], [9, 0, 1]]) # t = 0
X1_batch = np.array([[9, 8, 7], [0, 0, 0], [6, 5, 4], [3, 2, 1]]) # t = 1

X0 = tf.placeholder(tf.float32, [None, n_inputs])
X1 = tf.placeholder(tf.float32, [None, n_inputs])
```

▶ Manul implementation of an RNN

```
# make the dataset
n_inputs = 3
n_neurons = 5

X0_batch = np.array([[0, 1, 2], [3, 4, 5], [6, 7, 8], [9, 0, 1]]) # t = 0
X1_batch = np.array([[9, 8, 7], [0, 0, 0], [6, 5, 4], [3, 2, 1]]) # t = 1

X0 = tf.placeholder(tf.float32, [None, n_inputs])
X1 = tf.placeholder(tf.float32, [None, n_inputs])
```

```
# build the network
Wx = tf.Variable(tf.random_normal(shape=[n_inputs, n_neurons], dtype=tf.float32))
Wh = tf.Variable(tf.random_normal(shape=[n_neurons, n_neurons], dtype=tf.float32))
b = tf.Variable(tf.zeros([1, n_neurons], dtype=tf.float32))

h0 = tf.tanh(tf.matmul(X0, Wx) + b)
h1 = tf.tanh(tf.matmul(h0, Wh) + tf.matmul(X1, Wx) + b)
```

- Use `dynamic_rnn`

```
n_inputs = 3
n_neurons = 5
n_steps = 2

X_batch = np.array([
        # t = 0        t = 1
        [[0, 1, 2], [9, 8, 7]], # instance 1
        [[3, 4, 5], [0, 0, 0]], # instance 2
        [[6, 7, 8], [6, 5, 4]], # instance 3
        [[9, 0, 1], [3, 2, 1]], # instance 4
    ])

X = tf.placeholder(tf.float32, [None, n_steps, n_inputs])
```

▶ Use `dynamic_rnn`

```
n_inputs = 3
n_neurons = 5
n_steps = 2

X_batch = np.array([
        # t = 0       t = 1
        [[0, 1, 2], [9, 8, 7]], # instance 1
        [[3, 4, 5], [0, 0, 0]], # instance 2
        [[6, 7, 8], [6, 5, 4]], # instance 3
        [[9, 0, 1], [3, 2, 1]], # instance 4
    ])

X = tf.placeholder(tf.float32, [None, n_steps, n_inputs])
```

```
# build the network
basic_cell = tf.contrib.rnn.BasicRNNCell(num_units=n_neurons)
outputs, states = tf.nn.dynamic_rnn(basic_cell, X, dtype=tf.float32)
```

- Multi-layer RNN

```
layers = [tf.contrib.rnn.BasicRNNCell(num_units=n_neurons, activation=tf.nn.relu)
  for layer in range(n_layers)]

multi_layer_cell = tf.contrib.rnn.MultiRNNCell(layers)

outputs, states = tf.nn.dynamic_rnn(multi_layer_cell, X, dtype=tf.float32)

states_concat = tf.concat(axis=1, values=states)

logits = tf.layers.dense(states_concat, n_outputs)
```
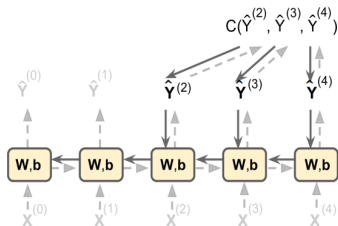
# Training RNNs

- To train an RNN, we should unroll it through time and then simply use regular backpropagation.
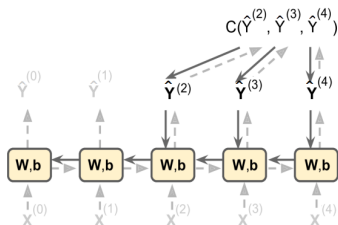
- This strategy is called backpropagation through time (BPTT).

▶ To train the model using BPTT, we go through the following steps:

▶ 1. Forward pass through the unrolled network (represented by the dashed arrows).

▶ 2. The cost function is $C(\hat{\mathbf{y}}^{\mathtt{tmin}}, \hat{\mathbf{y}}^{\mathtt{tmin+1}}, \cdots, \hat{\mathbf{y}}^{\mathtt{tmax}})$, where tmin and tmax are the first and last output time steps, not counting the ignored outputs.
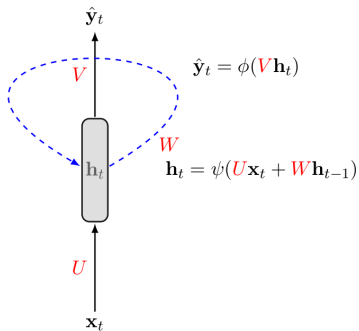
- 3. Propagate backward the gradients of that cost function through the unrolled network (represented by the solid arrows).

- 4. The model parameters are updated using the gradients computed during BPTT.

▶ The gradients flow backward through all the outputs used by the cost function, not just through the final output.

▶ For example, in the following figure:
  • The cost function is computed using the last three outputs, $\hat{\mathbf{y}}^{(2)}$, $\hat{\mathbf{y}}^{(3)}$, and $\hat{\mathbf{y}}^{(4)}$.
  • Gradients flow through these three outputs, but not through $\hat{\mathbf{y}}^{(0)}$ and $\hat{\mathbf{y}}^{(1)}$.

$$\hat{\mathbf{y}}_t = \phi(V\mathbf{h}_t)$$

$$\mathbf{h}_t = \psi(U\mathbf{x}_t + W\mathbf{h}_{t-1})$$

$$\mathbf{x}_1 \qquad \mathbf{x}_2 \qquad \mathbf{x}_3 \qquad \cdots \qquad \mathbf{x}_\tau$$

$$\mathbf{s}^{(t)} = \mathbf{u}^T\mathbf{x}^{(t)} + \mathbf{w}\mathbf{h}^{(t-1)}$$

$$\mathbf{h}^{(t)} = \tanh(\mathbf{s}^{(t)})$$

$$\mathbf{z}^{(t)} = \mathbf{v}\mathbf{h}^{(t)}$$

$$\hat{\mathbf{y}}^{(t)} = \texttt{softmax}(\mathbf{z}^{(t)})$$

$$\mathbf{J}^{(t)} = \texttt{cross\_entropy}(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = -\sum \mathbf{y}^{(t)}\log\hat{\mathbf{y}}^{(t)}$$
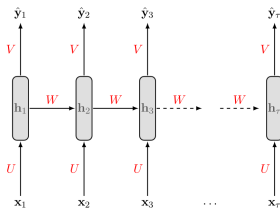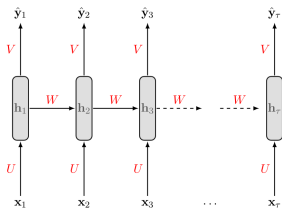
$$J^{(t)} = \texttt{cross\_entropy}(y^{(t)}, \hat{y}^{(t)}) = -\sum y^{(t)} \log \hat{y}^{(t)}$$

▶ We treat the full sequence as one training example.

$$J^{(t)} = \texttt{cross\_entropy}(y^{(t)}, \hat{y}^{(t)}) = -\sum y^{(t)} \log \hat{y}^{(t)}$$

▶ We treat the full sequence as one training example.

▶ The total error E is just the sum of the errors at each time step.

▶ E.g., $E = J^{(1)} + J^{(2)} + \cdots + J^{(t)}$

▶ $J^{(t)}$ is the total cost, so we can say that a 1-unit increase in $v$, $w$ or $u$ will impact each of $J^{(1)}$, $J^{(2)}$, until $J^{(t)}$ individually.

▶ $J^{(t)}$ is the total cost, so we can say that a 1-unit increase in $v$, $w$ or $u$ will impact each of $J^{(1)}$, $J^{(2)}$, until $J^{(t)}$ individually.

▶ The gradient is equal to the sum of the respective gradients at each time step $t$.

▶ $J^{(t)}$ is the total cost, so we can say that a 1-unit increase in $v$, $w$ or $u$ will impact each of $J^{(1)}$, $J^{(2)}$, until $J^{(t)}$ individually.

▶ The gradient is equal to the sum of the respective gradients at each time step $t$.

▶ For example if $t = 3$ we have: $E = J^{(1)} + J^{(2)} + J^{(3)}$

▶ $J^{(t)}$ is the total cost, so we can say that a 1-unit increase in v, w or u will impact each of $J^{(1)}$, $J^{(2)}$, until $J^{(t)}$ individually.

▶ The gradient is equal to the sum of the respective gradients at each time step t.

▶ For example if t = 3 we have: $E = J^{(1)} + J^{(2)} + J^{(3)}$

$$\frac{\partial E}{\partial v} = \sum_t \frac{\partial J^{(t)}}{\partial v} = \frac{\partial J^{(3)}}{\partial v} + \frac{\partial J^{(2)}}{\partial v} + \frac{\partial J^{(1)}}{\partial v}$$

▶ $J^{(t)}$ is the total cost, so we can say that a 1-unit increase in $v$, $w$ or $u$ will impact each of $J^{(1)}$, $J^{(2)}$, until $J^{(t)}$ individually.

▶ The gradient is equal to the sum of the respective gradients at each time step $t$.

▶ For example if $t = 3$ we have: $E = J^{(1)} + J^{(2)} + J^{(3)}$

$$\frac{\partial E}{\partial v} = \sum_t \frac{\partial J^{(t)}}{\partial v} = \frac{\partial J^{(3)}}{\partial v} + \frac{\partial J^{(2)}}{\partial v} + \frac{\partial J^{(1)}}{\partial v}$$

$$\frac{\partial E}{\partial w} = \sum_t \frac{\partial J^{(t)}}{\partial w} = \frac{\partial J^{(3)}}{\partial w} + \frac{\partial J^{(2)}}{\partial w} + \frac{\partial J^{(1)}}{\partial w}$$

▶ $J^{(t)}$ is the total cost, so we can say that a 1-unit increase in `v`, `w` or `u` will impact each of $J^{(1)}$, $J^{(2)}$, until $J^{(t)}$ individually.

▶ The gradient is equal to the sum of the respective gradients at each time step `t`.

▶ For example if `t = 3` we have: $E = J^{(1)} + J^{(2)} + J^{(3)}$

$$\frac{\partial E}{\partial v} = \sum_t \frac{\partial J^{(t)}}{\partial v} = \frac{\partial J^{(3)}}{\partial v} + \frac{\partial J^{(2)}}{\partial v} + \frac{\partial J^{(1)}}{\partial v}$$
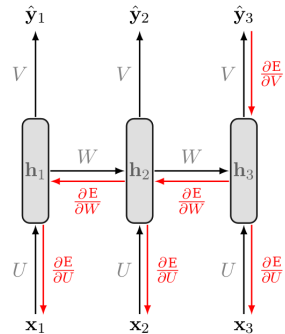
$$\frac{\partial E}{\partial w} = \sum_t \frac{\partial J^{(t)}}{\partial w} = \frac{\partial J^{(3)}}{\partial w} + \frac{\partial J^{(2)}}{\partial w} + \frac{\partial J^{(1)}}{\partial w}$$

$$\frac{\partial E}{\partial u} = \sum_t \frac{\partial J^{(3)}}{\partial u} = \frac{\partial J^{(3)}}{\partial u} + \frac{\partial J^{(2)}}{\partial u} + \frac{\partial J^{(1)}}{\partial u}$$

- Let's start with $\frac{\partial \mathbf{E}}{\partial \mathbf{v}}$.

- A change in $\mathbf{v}$ will only impact $\mathbf{J}^{(3)}$ at time $\mathbf{t} = 3$, because it plays no role in computing the value of anything other than $\mathbf{z}^{(3)}$.

$$\frac{\partial \mathbf{E}}{\partial \mathbf{v}} = \sum_{\mathbf{t}} \frac{\partial \mathbf{J}^{(\mathbf{t})}}{\partial \mathbf{v}} = \frac{\partial \mathbf{J}^{(3)}}{\partial \mathbf{v}} + \frac{\partial \mathbf{J}^{(2)}}{\partial \mathbf{v}} + \frac{\partial \mathbf{J}^{(1)}}{\partial \mathbf{v}}$$

- Let's start with $\frac{\partial E}{\partial v}$.

- A change in $v$ will only impact $J^{(3)}$ at time $t = 3$, because it plays no role in computing the value of anything other than $z^{(3)}$.

$$\frac{\partial E}{\partial v} = \sum_t \frac{\partial J^{(t)}}{\partial v} = \frac{\partial J^{(3)}}{\partial v} + \frac{\partial J^{(2)}}{\partial v} + \frac{\partial J^{(1)}}{\partial v}$$
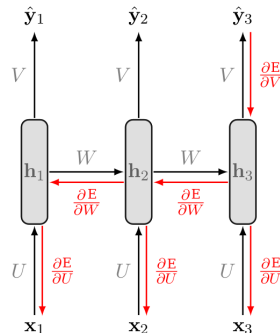
$$\frac{\partial J^{(3)}}{\partial v} = \frac{\partial J^{(3)}}{\partial \hat{y}^{(3)}} \frac{\partial \hat{y}^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial v}$$
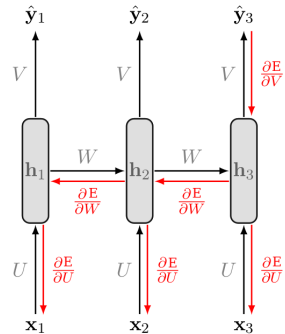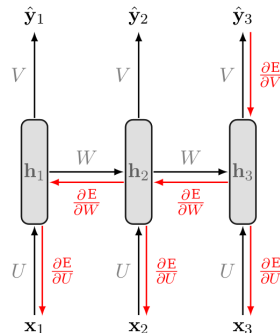
▶ Let's start with $\frac{\partial E}{\partial v}$.

▶ A change in $v$ will only impact $J^{(3)}$ at time $t = 3$, because it plays no role in computing the value of anything other than $z^{(3)}$.

$$\frac{\partial E}{\partial v} = \sum_t \frac{\partial J^{(t)}}{\partial v} = \frac{\partial J^{(3)}}{\partial v} + \frac{\partial J^{(2)}}{\partial v} + \frac{\partial J^{(1)}}{\partial v}$$

$$\frac{\partial J^{(3)}}{\partial v} = \frac{\partial J^{(3)}}{\partial \hat{y}^{(3)}} \frac{\partial \hat{y}^{(3)}}{\partial z^{(3)}} \frac{\partial z^{(3)}}{\partial v}$$
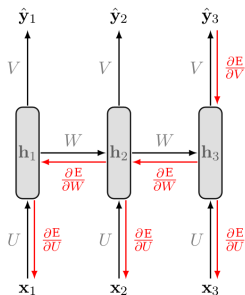
$$\frac{\partial J^{(2)}}{\partial v} = \frac{\partial J^{(2)}}{\partial \hat{y}^{(2)}} \frac{\partial \hat{y}^{(2)}}{\partial z^{(2)}} \frac{\partial z^{(2)}}{\partial v}$$

► Let's start with $\frac{\partial \mathbf{E}}{\partial \mathbf{v}}$.

► A change in $\mathbf{v}$ will only impact $\mathtt{J}^{(3)}$ at time $\mathtt{t} = 3$, because it plays no role in computing the value of anything other than $\mathbf{z}^{(3)}$.

$$\frac{\partial \mathbf{E}}{\partial \mathbf{v}} = \sum_{\mathtt{t}} \frac{\partial \mathtt{J}^{(\mathtt{t})}}{\partial \mathbf{v}} = \frac{\partial \mathtt{J}^{(3)}}{\partial \mathbf{v}} + \frac{\partial \mathtt{J}^{(2)}}{\partial \mathbf{v}} + \frac{\partial \mathtt{J}^{(1)}}{\partial \mathbf{v}}$$

$$\frac{\partial \mathtt{J}^{(3)}}{\partial \mathbf{v}} = \frac{\partial \mathtt{J}^{(3)}}{\partial \hat{\mathbf{y}}^{(3)}} \frac{\partial \hat{\mathbf{y}}^{(3)}}{\partial \mathbf{z}^{(3)}} \frac{\partial \mathbf{z}^{(3)}}{\partial \mathbf{v}}$$

$$\frac{\partial \mathtt{J}^{(2)}}{\partial \mathbf{v}} = \frac{\partial \mathtt{J}^{(2)}}{\partial \hat{\mathbf{y}}^{(2)}} \frac{\partial \hat{\mathbf{y}}^{(2)}}{\partial \mathbf{z}^{(2)}} \frac{\partial \mathbf{z}^{(2)}}{\partial \mathbf{v}}$$

$$\frac{\partial \mathtt{J}^{(1)}}{\partial \mathbf{v}} = \frac{\partial \mathtt{J}^{(1)}}{\partial \hat{\mathbf{y}}^{(1)}} \frac{\partial \hat{\mathbf{y}}^{(1)}}{\partial \mathbf{z}^{(1)}} \frac{\partial \mathbf{z}^{(1)}}{\partial \mathbf{v}}$$

▶ Let's compute the derivatives of $\frac{\partial J}{\partial w}$ and $\frac{\partial J}{\partial u}$, which are computed the same.

▶ A change in w at t = 3 will impact our cost J in 3 separate ways:

1. When computing the value of $h^{(1)}$.
2. When computing the value of $h^{(2)}$, which depends on $h^{(1)}$.
3. When computing the value of $h^{(3)}$, which depends on $h^{(2)}$, which depends on $h^{(1)}$.
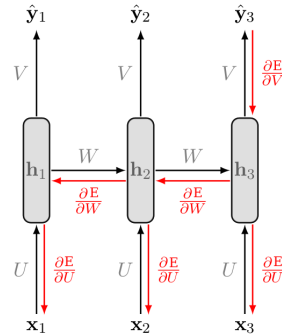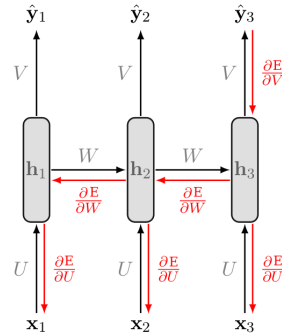
▶ we compute our individual gradients as:

$$\sum_{t} \frac{\partial J^{(t)}}{\partial w} = \frac{\partial J^{(3)}}{\partial w} + \frac{\partial J^{(2)}}{\partial w} + \frac{\partial J^{(1)}}{\partial w}$$

$$\frac{\partial J^{(1)}}{\partial w} = \frac{\partial J^{(1)}}{\partial \hat{y}^{(1)}} \frac{\partial \hat{y}^{(1)}}{\partial z^{(1)}} \frac{\partial z^{(1)}}{\partial h^{(1)}} \frac{\partial h^{(1)}}{\partial s^{(1)}} \frac{\partial s^{(1)}}{\partial w}$$

▶ we compute our individual gradients as:

$$\sum_{\mathtt{t}} \frac{\partial \mathtt{J}^{(\mathtt{t})}}{\partial \mathtt{w}} = \frac{\partial \mathtt{J}^{(3)}}{\partial \mathtt{w}} + \frac{\partial \mathtt{J}^{(2)}}{\partial \mathtt{w}} + \frac{\partial \mathtt{J}^{(1)}}{\partial \mathtt{w}}$$

$$\frac{\partial \mathtt{J}^{(2)}}{\partial \mathtt{w}} = \frac{\partial \mathtt{J}^{(2)}}{\partial \hat{\mathtt{y}}^{(2)}} \frac{\partial \hat{\mathtt{y}}^{(2)}}{\partial \mathtt{z}^{(2)}} \frac{\partial \mathtt{z}^{(2)}}{\partial \mathtt{h}^{(2)}} \frac{\partial \mathtt{h}^{(2)}}{\partial \mathtt{s}^{(2)}} \frac{\partial \mathtt{s}^{(2)}}{\partial \mathtt{w}} +$$

$$\frac{\partial \mathtt{J}^{(2)}}{\partial \hat{\mathtt{y}}^{(2)}} \frac{\partial \hat{\mathtt{y}}^{(2)}}{\partial \mathtt{z}^{(2)}} \frac{\partial \mathtt{z}^{(2)}}{\partial \mathtt{h}^{(2)}} \frac{\partial \mathtt{h}^{(2)}}{\partial \mathtt{s}^{(2)}} \frac{\partial \mathtt{s}^{(2)}}{\partial \mathtt{h}^{(1)}} \frac{\partial \mathtt{h}^{(1)}}{\partial \mathtt{s}^{(1)}} \frac{\partial \mathtt{s}^{(1)}}{\partial \mathtt{w}}$$
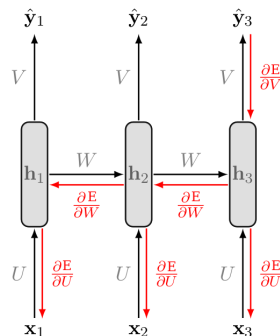
▶ we compute our individual gradients as:

$$\sum_t \frac{\partial \mathtt{J}^{(t)}}{\partial \mathtt{w}} = \frac{\partial \mathtt{J}^{(3)}}{\partial \mathtt{w}} + \frac{\partial \mathtt{J}^{(2)}}{\partial \mathtt{w}} + \frac{\partial \mathtt{J}^{(1)}}{\partial \mathtt{w}}$$

$$\frac{\partial \mathtt{J}^{(3)}}{\partial \mathtt{w}} = \frac{\partial \mathtt{J}^{(3)}}{\partial \hat{\mathtt{y}}^{(3)}} \frac{\partial \hat{\mathtt{y}}^{(3)}}{\partial \mathtt{z}^{(3)}} \frac{\partial \mathtt{z}^{(3)}}{\partial \mathtt{h}^{(3)}} \frac{\partial \mathtt{h}^{(3)}}{\partial \mathtt{s}^{(3)}} \frac{\partial \mathtt{s}^{(3)}}{\partial \mathtt{w}} +$$
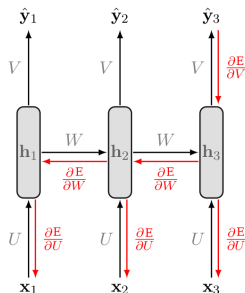
$$\frac{\partial \mathtt{J}^{(3)}}{\partial \hat{\mathtt{y}}^{(3)}} \frac{\partial \hat{\mathtt{y}}^{(3)}}{\partial \mathtt{z}^{(3)}} \frac{\partial \mathtt{z}^{(3)}}{\partial \mathtt{h}^{(3)}} \frac{\partial \mathtt{h}^{(3)}}{\partial \mathtt{s}^{(3)}} \frac{\partial \mathtt{s}^{(3)}}{\partial \mathtt{h}^{(2)}} \frac{\partial \mathtt{s}^{(t)}}{\partial \mathtt{w}} +$$

$$\frac{\partial \mathtt{J}^{(3)}}{\partial \hat{\mathtt{y}}^{(3)}} \frac{\partial \hat{\mathtt{y}}^{(3)}}{\partial \mathtt{z}^{(3)}} \frac{\partial \mathtt{z}^{(3)}}{\partial \mathtt{h}^{(3)}} \frac{\partial \mathtt{h}^{(3)}}{\partial \mathtt{s}^{(3)}} \frac{\partial \mathtt{s}^{(3)}}{\partial \mathtt{h}^{(2)}} \frac{\partial \mathtt{h}^{(2)}}{\partial \mathtt{s}^{(2)}} \frac{\partial \mathtt{s}^{(2)}}{\partial \mathtt{h}^{(1)}} \frac{\partial \mathtt{h}^{(1)}}{\partial \mathtt{s}^{(1)}} \frac{\partial \mathtt{s}^{(1)}}{\partial \mathtt{w}}$$

▶ More generally, a change in `w` will impact our cost `J⁽ᵗ⁾` on `t` separate occasions.

$$\frac{\partial \mathtt{J}^{(\mathtt{t})}}{\partial \mathtt{w}} = \sum_{k=1}^{t} \frac{\partial \mathtt{J}^{(\mathtt{t})}}{\partial \hat{\mathtt{y}}^{(\mathtt{t})}} \frac{\partial \hat{\mathtt{y}}^{(\mathtt{t})}}{\partial \mathtt{z}^{(\mathtt{t})}} \frac{\partial \hat{\mathtt{z}}^{(\mathtt{t})}}{\partial \mathtt{h}^{(\mathtt{t})}} \left( \prod_{j=k+1}^{t} \frac{\partial \mathtt{h}^{(\mathtt{j})}}{\partial \mathtt{s}^{(\mathtt{j})}} \frac{\partial \mathtt{s}^{(\mathtt{j})}}{\partial \mathtt{h}^{(\mathtt{j}-1)}} \right) \frac{\partial \mathtt{h}^{(\mathtt{k})}}{\partial \mathtt{s}^{(\mathtt{k})}} \frac{\partial \mathtt{s}^{(\mathtt{k})}}{\partial \mathtt{w}}$$

# RNN Design Patterns

▶ Sequence-to-vector network: takes a sequence of inputs, and ignore all outputs except for the last one.

- **Sequence-to-vector** network: takes a sequence of inputs, and ignore all outputs except for the last one.

- E.g., you could feed the network a sequence of words corresponding to a movie review, and the network would output a sentiment score.
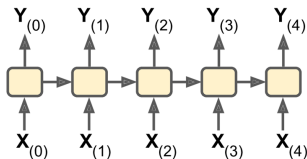
▶ Vector-to-sequence network: takes a single input at the first time step, and let it output a sequence.

- Vector-to-sequence network: takes a single input at the first time step, and let it output a sequence.
- E.g., the input could be an image, and the output could be a caption for that image.

- Sequence-to-sequence network: takes a sequence of inputs and produce a sequence of outputs.

- **Sequence-to-sequence** network: takes a sequence of inputs and produce a sequence of outputs.

- Useful for predicting time series such as stock prices: you feed it the prices over the last N days, and it must output the prices shifted by one day into the future.

- Here, both input sequences and output sequences have the same length.

- Encoder-decoder network: a sequence-to-vector network (encoder), followed by a vector-to-sequence network (decoder).

- Encoder-decoder network: a sequence-to-vector network (encoder), followed by a vector-to-sequence network (decoder).

- E.g., translating a sentence from one language to another.

- You would feed the network a sentence in one language, the encoder would convert this sentence into a single vector representation, and then the decoder would decode this vector into a sentence in another language.

# LSTM

▶ Sometimes we only need to look at recent information to perform the present task.
  • E.g., predicting the next word based on the previous ones.

# RNN Problems

- Sometimes we only need to look at recent information to perform the present task.
  - E.g., predicting the next word based on the previous ones.

- In such cases, where the gap between the relevant information and the place that it's needed is small, RNNs can learn to use the past information.

- Sometimes we only need to look at recent information to perform the present task.
  - E.g., predicting the next word based on the previous ones.

- In such cases, where the gap between the relevant information and the place that it's needed is small, RNNs can learn to use the past information.

- But, as that gap grows, RNNs become unable to learn to connect the information.

- RNNs may suffer from the vanishing/exploding gradients problem.

# RNN Problems

▶ Sometimes we only need to look at recent information to perform the present task.
  • E.g., predicting the next word based on the previous ones.

▶ In such cases, where the gap between the relevant information and the place that it's needed is small, RNNs can learn to use the past information.

▶ But, as that gap grows, RNNs become unable to learn to connect the information.

▶ RNNs may suffer from the vanishing/exploding gradients problem.

▶ To solve these problem, long short-term memory (LSTM) have been introduced.

# RNN Problems

- Sometimes we only need to look at recent information to perform the present task.
  - E.g., predicting the next word based on the previous ones.

- In such cases, where the gap between the relevant information and the place that it's needed is small, RNNs can learn to use the past information.

- But, as that gap grows, RNNs become unable to learn to connect the information.

- RNNs may suffer from the vanishing/exploding gradients problem.

- To solve these problem, long short-term memory (LSTM) have been introduced.

- In LSTM, the network can learn what to store and what to throw away.

- Without looking inside the box, the LSTM cell looks exactly like a basic cell.

# RNN Basic Cell vs. LSTM

- Without looking inside the box, the LSTM cell looks exactly like a basic cell.

- The repeating module in a standard RNN contains a single layer.

▶ Without looking inside the box, the LSTM cell looks exactly like a basic cell.

▶ The repeating module in a standard RNN contains a single layer.



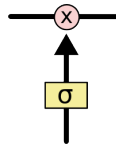▶ The repeating module in an LSTM contains four interacting layers.

# LSTM (1/2)

- In LSTM state is split in two vectors:
  1. $h^{(t)}$ (h stands for hidden): the short-term state
  2. $c^{(t)}$ (c stands for cell): the long-term state

# LSTM (2/2)

- ▶ The cell state (long-term state), the horizontal line on the top of the diagram.
- ▶ The LSTM can remove/add information to the cell state, regulated by three gates.
  - Forget gate, input gate and output gate

▶ **Step one**: decides what information we are going to throw away from the cell state.

$$\mathbf{f}^{(t)} = \sigma(\mathbf{u}_f^\mathsf{T} \mathbf{x}^{(t)} + \mathbf{w}_f \mathbf{h}^{(t-1)})$$

▶ **Step one**: decides what information we are going to throw away from the cell state.

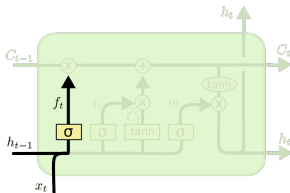▶ This decision is made by a sigmoid layer, called the forget gate layer.

$$\mathbf{f^{(t)}} = \sigma(\mathbf{u_f^\top x^{(t)}} + \mathbf{w_f h^{(t-1)}})$$

▶ **Step one**: decides what information we are going to throw away from the cell state.

▶ This decision is made by a sigmoid layer, called the forget gate layer.

▶ It looks at $h^{(t-1)}$ and $\mathbf{x}^{(t)}$, and outputs a number between 0 and 1 for each number in the cell state $c^{(t-1)}$.

  • 1 represents completely keep this, and 0 represents completely get rid of this.

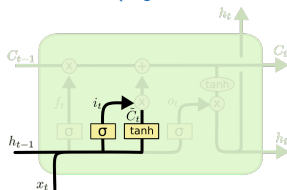$$f^{(t)} = \sigma(\mathbf{u}_f^\mathsf{T}\mathbf{x}^{(t)} + \mathbf{w}_f h^{(t-1)})$$

▶ Second step: decides what new information we are going to store in the cell state. This has two parts:

$$i^{(t)} = \sigma(\mathbf{u}_i^{\mathsf{T}}\mathbf{x}^{(t)} + \mathtt{w}_i h^{(t-1)})$$

$$\tilde{c}^{(t)} = \tanh(\mathbf{u}_{\tilde{c}}^{\mathsf{T}}\mathbf{x}^{(t)} + \mathtt{w}_{\tilde{c}} h^{(t-1)})$$

▶ Second step: decides what new information we are going to store in the cell state. This has two parts:

▶ 1. A sigmoid layer, called the input gate layer, decides which values we will update.

$$i^{(t)} = \sigma(\mathbf{u}_i^\intercal \mathbf{x}^{(t)} + \mathtt{w}_i h^{(t-1)})$$

$$\tilde{c}^{(t)} = \tanh(\mathbf{u}_{\tilde{c}}^\intercal \mathbf{x}^{(t)} + \mathtt{w}_{\tilde{c}} h^{(t-1)})$$

- Second step: decides what new information we are going to store in the cell state. This has two parts:

- 1. A sigmoid layer, called the input gate layer, decides which values we will update.

- 2. A tanh layer creates a vector of new candidate values that could be added to the state.

$$i^{(t)} = \sigma(\mathbf{u}_i^\intercal \mathbf{x}^{(t)} + \mathbf{w}_i h^{(t-1)})$$

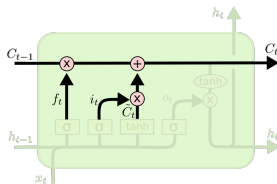$$\tilde{c}^{(t)} = \tanh(\mathbf{u}_{\tilde{c}}^\intercal \mathbf{x}^{(t)} + \mathbf{w}_{\tilde{c}} h^{(t-1)})$$

▶ Third step: updates the old cell state $c^{(t-1)}$, into the new cell state $c^{(t)}$.

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

▶ Third step: updates the old cell state $c^{(t-1)}$, into the new cell state $c^{(t)}$.

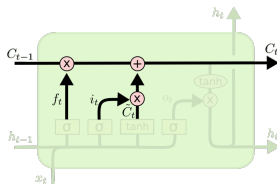▶ We multiply the old state by $f^{(t)}$, forgetting the things we decided to forget earlier.

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

- ▶ Third step: updates the old cell state $c^{(t-1)}$, into the new cell state $c^{(t)}$.
- ▶ We multiply the old state by $f^{(t)}$, forgetting the things we decided to forget earlier.
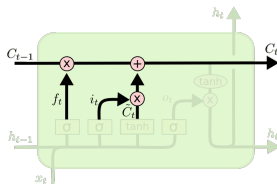- ▶ Then we add it $i^{(t)} \otimes \tilde{c}^{(t)}$.

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

- **Third step**: updates the old cell state $c^{(t-1)}$, into the new cell state $c^{(t)}$.

- We multiply the old state by $f^{(t)}$, forgetting the things we decided to forget earlier.

- Then we add it $i^{(t)} \otimes \tilde{c}^{(t)}$.

- This is the new candidate values, scaled by how much we decided to update each state value.

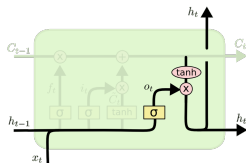$$c^{(t)} = f^{(t)} \otimes c^{(t-1)} + i^{(t)} \otimes \tilde{c}^{(t)}$$

► Fourth step: decides about the output.

$$o^{(t)} = \sigma(\mathbf{u}_o^\mathsf{T}\mathbf{x}^{(t)} + \mathrm{w}_o h^{(t-1)})$$
$$\hat{y}^{(t)} = h^{(t)} = o^{(t)} \otimes \tanh(c^{(t)})$$

- **Fourth step**: decides about the output.
- First, runs a sigmoid layer that decides what parts of the cell state we are going to output.

$$o^{(t)} = \sigma(\mathbf{u}_o^{\mathsf{T}}\mathbf{x}^{(t)} + w_o h^{(t-1)})$$
$$\hat{y}^{(t)} = h^{(t)} = o^{(t)} \otimes \tanh(c^{(t)})$$

▶ **Fourth step**: decides about the output.

▶ First, runs a sigmoid layer that decides what parts of the cell state we are going to output.

▶ Then, puts the cell state through tanh and multiplies it by the output of the sigmoid gate (output gate), so that it only outputs the parts it decided to.

$$o^{(t)} = \sigma(\mathbf{u}_o^{\intercal} \mathbf{x}^{(t)} + w_o h^{(t-1)})$$

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

▶ Multi-layer LSTM

```
lstm_cells = [tf.contrib.rnn.BasicLSTMCell(num_units=n_neurons) for layer in range(n_layers)]

multi_cell = tf.contrib.rnn.MultiRNNCell(lstm_cells)

outputs, states = tf.nn.dynamic_rnn(multi_cell, X, dtype=tf.float32)

top_layer_h_state = states[-1][1]

logits = tf.layers.dense(top_layer_h_state, n_outputs)
```

# Summary

# Summary

- RNN

- Unfolding the network

- Three weights

- Backpropagation through time

- RNN design patterns

- LSTM

# Reference

- Ian Goodfellow et al., Deep Learning (Ch. 10)

- Aurélien Géron, Hands-On Machine Learning (Ch. 14)

- Understanding LSTM Networks
  `http://colah.github.io/posts/2015-08-Understanding-LSTMs`

- CS224d: Deep Learning for Natural Language Processing
  `http://cs224d.stanford.edu`

Questions?