

Deep Learning: Models for Sequence Data (RNN and LSTM) and Autoencoders

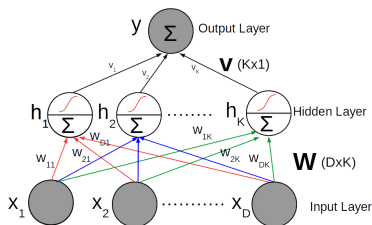
Piyush Rai

Machine Learning (CS771A)

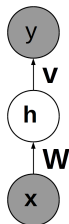
Nov 4, 2016

Recap: Feedforward Neural Network

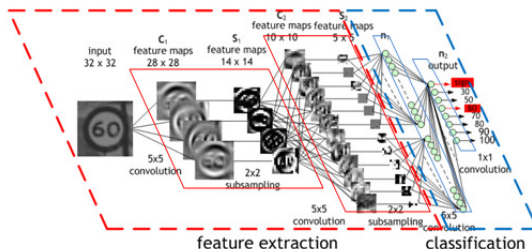
- Consists of an input layer, one or more hidden layers, and an output layer



- A “macro” view of the above (note: $\mathbf{x} = [x_1, \dots, x_D]$, $\mathbf{h} = [h_1, \dots, h_K]$)



Recap: Convolutional Neural Network

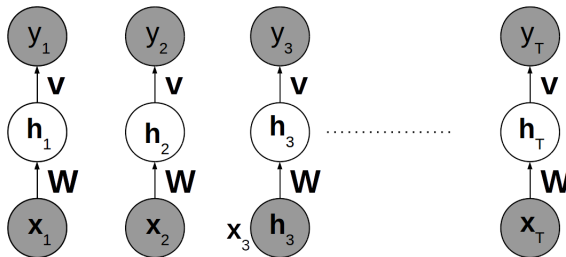


- Special type of feedforward neural nets (local connectivity + weight sharing)
- Each layer uses a set of “filters” (basically, weights to be learned) which can detect specific features. Filters are like basis/dictionary (PCA analogy)
- Each filter is convolved over entire input to produce a feature map
- Nonlinearity and pooling and applied after each convolution layer
- Last layer (one that connects to outputs) is fully connected

Deep Neural Networks for Modeling Sequence Data

Limitation of Feedforward Neural Nets

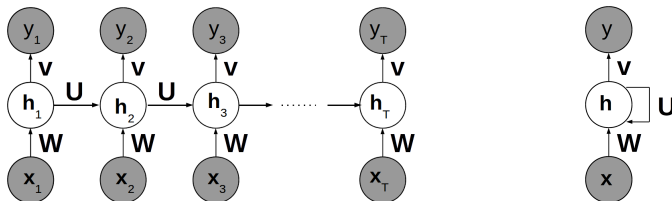
- FFNN can't take into account the sequential structure in the data
- For a sequence of observations $\mathbf{x}_1, \dots, \mathbf{x}_T$, their corresponding hidden units (states) $\mathbf{h}_1, \dots, \mathbf{h}_T$ are assumed independent of each other



- Not ideal for sequential data, e.g., sentence/paragraph/document (sequence of words), video (sequence of frames), etc.

Recurrent Neural Nets (RNN)

- Hidden state at each step depends on the hidden state of the previous



- Each hidden state is typically defined as

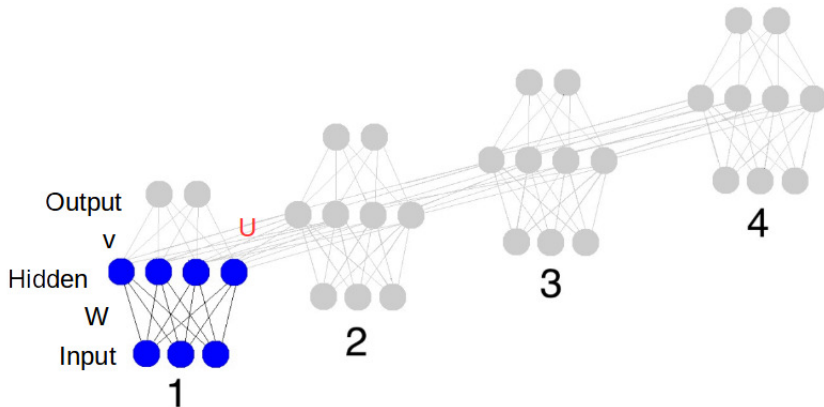
$$h_t = f(Wx_t + Uh_{t-1})$$

where U is like a transition matrix and f is some nonlin. fn. (e.g., tanh)

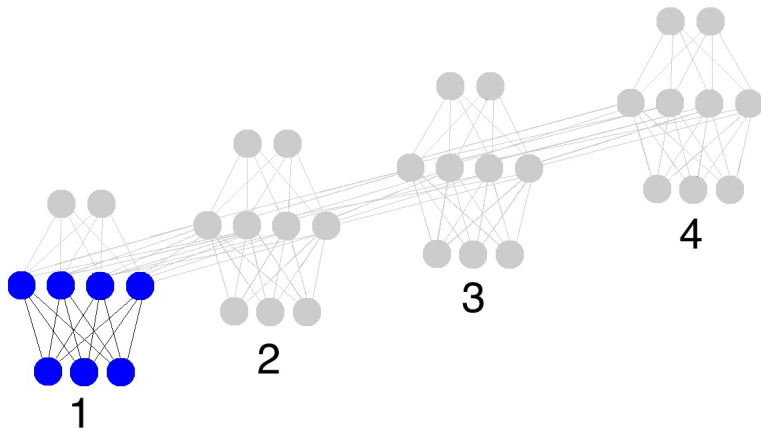
- Now h_t acts as a memory. Helps us remember what happened up to step t
- Note: Unlike sequence data models such as [HMM](#) where each state is discrete, RNN states are continuous-valued (in that sense, RNNs are similar to Linear-Gaussian models like [Kalman Filters](#) which have continuous states)
- RNNs can also be extended to have more than one hidden layer

Recurrent Neural Nets (RNN)

- A more “micro” view of RNN (the transition matrix \mathbf{U} connects the hidden states across observations, propagating information along the sequence)



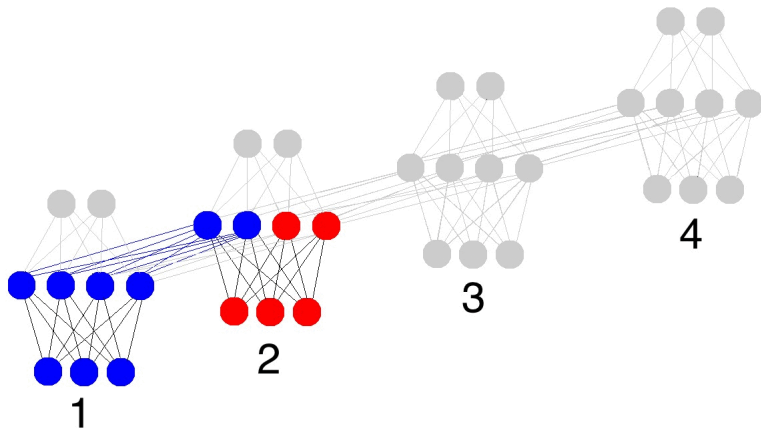
RNN in Action..



MakeAGIF.com

Pic courtesy: <https://iamtrask.github.io/>

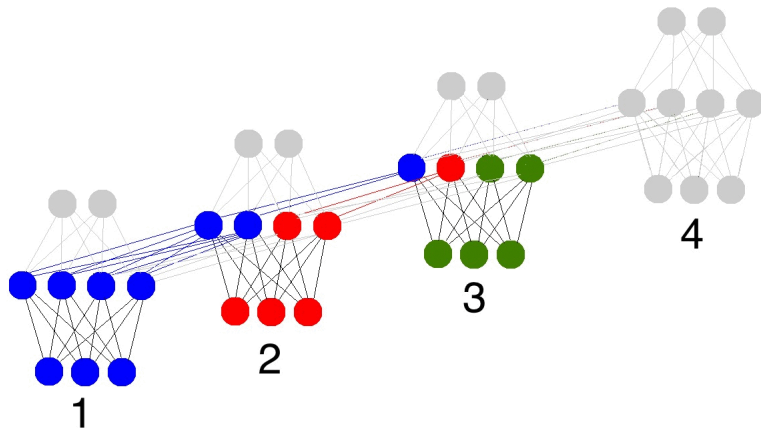
RNN in Action..



MakeAGIF.com

Pic courtesy: <https://iamtrask.github.io/>

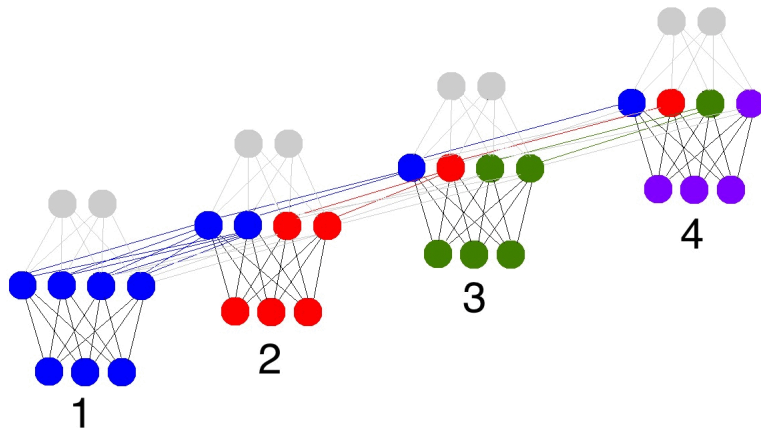
RNN in Action..



MakeAGIF.com

Pic courtesy: <https://iamtrask.github.io/>

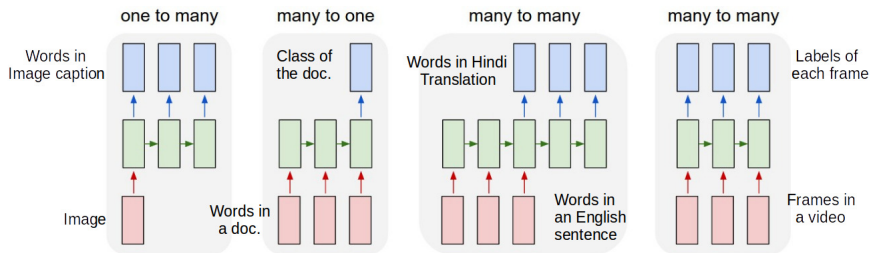
RNN in Action..



MakeAGIF.com

Pic courtesy: <https://iamtrask.github.io/>

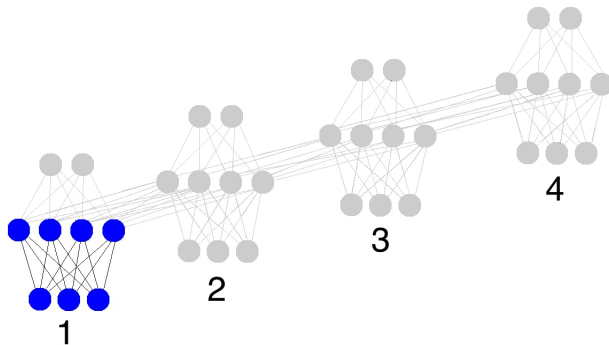
RNN: Applications



- RNNs are widely applicable and are also very flexible. E.g.,
 - Input, output, or both, can be sequences (possibly of different lengths)
 - Different inputs (and different outputs) need not be of the same length
 - Regardless of the length of the input sequence, RNN will learn a fixed size embedding for the input sequence

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

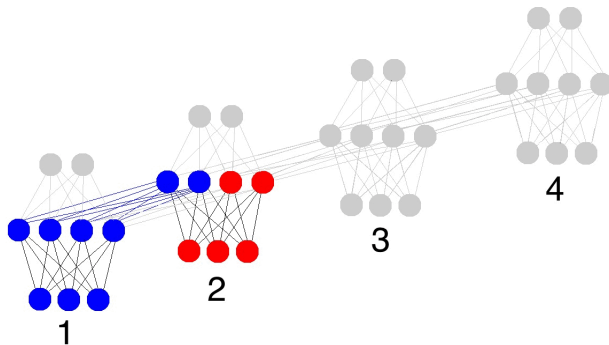


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

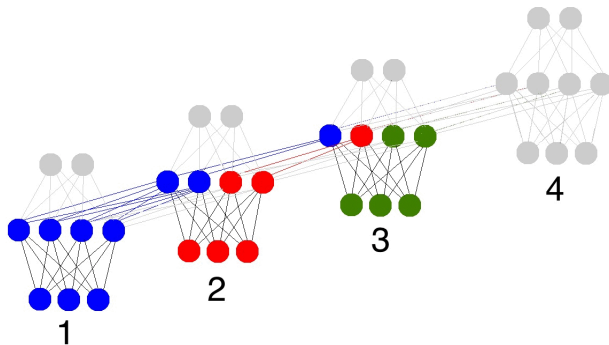


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

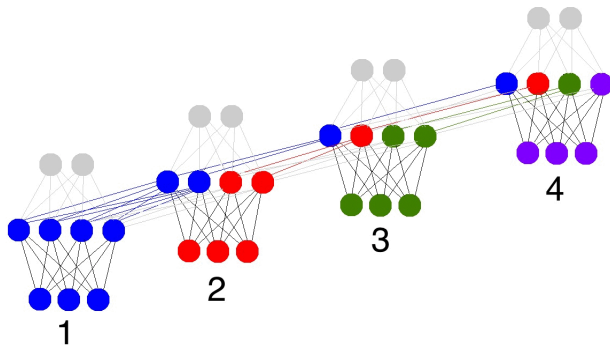


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

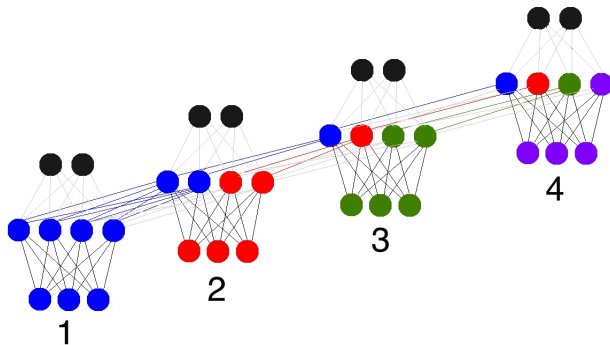


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

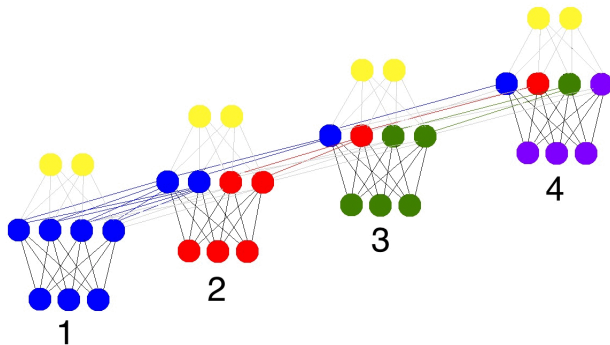


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

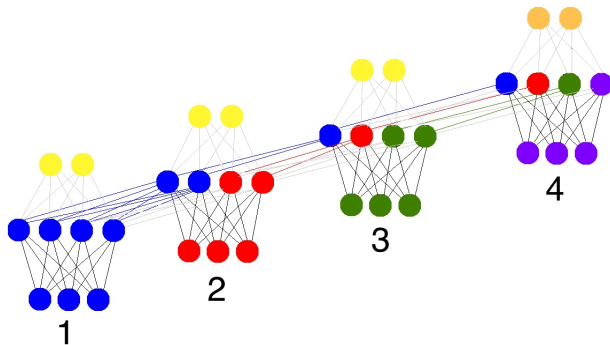


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

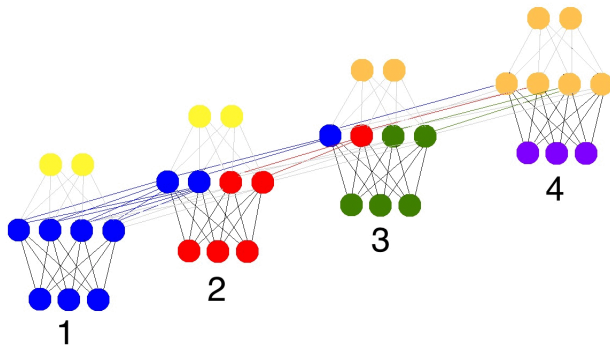


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

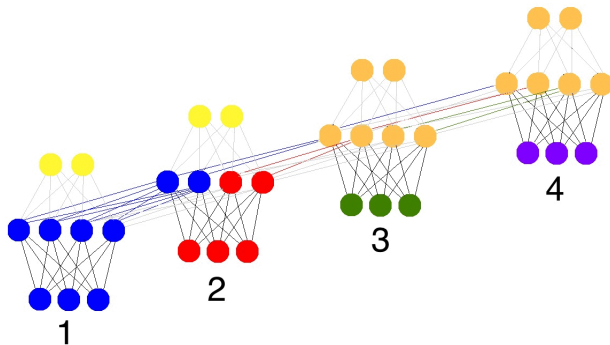


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

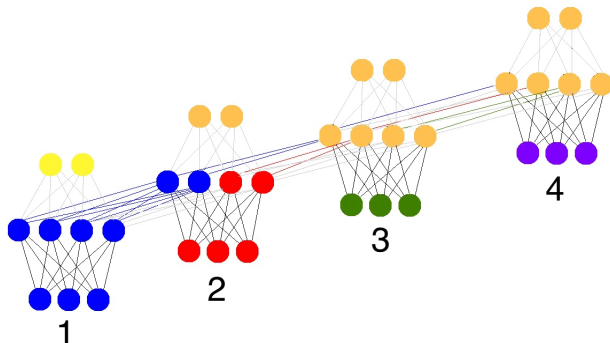


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

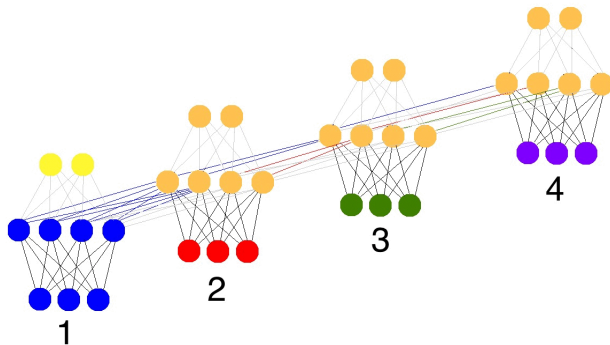


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

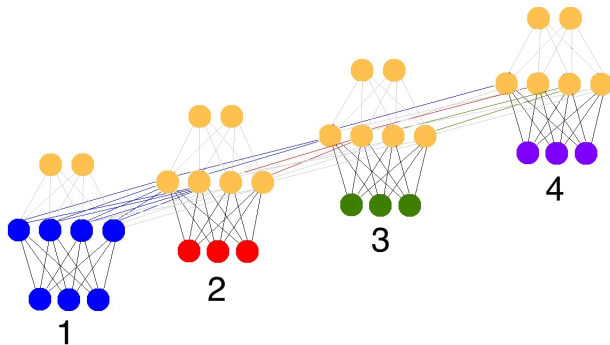


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets

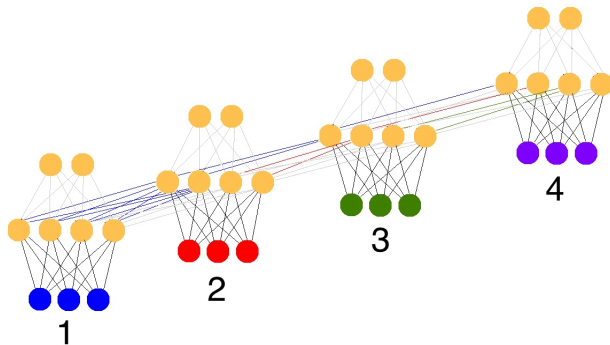


- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

Training RNN

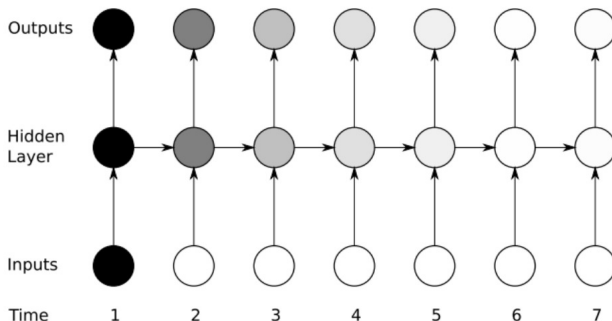
- Trained using **Backpropagation Through Time** (forward propagate from step 1 to end, and then backward propagate from end to step 1)
- Think of the time-dimension as another hidden layer and then it is just like standard backpropagation for feedforward neural nets



- Black: Prediction, Yellow: Error, Orange: Gradients

Pic courtesy: <https://iamtrask.github.io/>

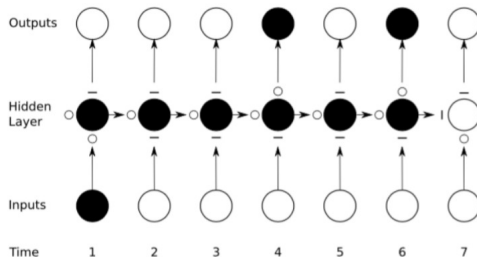
RNN: Vanishing/Exploding Gradients Problem



- Sensitivity of hidden states and outputs on a given input becomes weaker as we move away from it along the sequence (weak memory)
- New inputs “overwrite” the activations of previous hidden states
- Repeated multiplications can cause the gradients to vanish or explode

Capturing Long-Range Dependencies

- Idea: Augment the hidden states with **gates** (with parameters to be learned)
- These gates can help us remember and forget information “selectively”



- The hidden states have 3 type of gates
 - Input (bottom), Forget (left), Output (top)
- Open gate denoted by 'o', closed gate denoted by '-'
- **LSTM** (Hochreiter and Schmidhuber, mid-90s): **Long Short-Term Memory** is one such idea

Long Short-Term Memory (LSTM)

- Essentially an RNN, except that the hidden states are computed differently
- Recall that RNN computes the hidden states as $\mathbf{h}_t = \tanh(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1})$
- For RNN: State update is multiplicative (weak memory and gradient issues)
- In contrast, LSTM maintains a “context” C_t and computes hidden states as

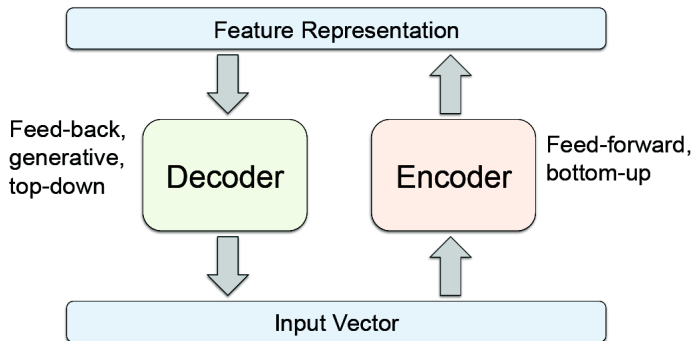
$$\begin{aligned}\hat{C}_t &= \tanh(\mathbf{W}^c \mathbf{x}_t + \mathbf{U}^c \mathbf{h}_{t-1}) && \text{ (“local” context, based on only the previous state)} \\ i_t &= \sigma(\mathbf{W}^i \mathbf{x}_t + \mathbf{U}^i \mathbf{h}_{t-1}) && \text{ (how much to take in the local context)} \\ f_t &= \sigma(\mathbf{W}^f \mathbf{x}_t + \mathbf{U}^f \mathbf{h}_{t-1}) && \text{ (how much to keep/forget the previous context)} \\ o_t &= \sigma(\mathbf{W}^o \mathbf{x}_t + \mathbf{U}^o \mathbf{h}_{t-1}) && \text{ (how much to output)} \\ C_t &= C_{t-1} \odot f_t + \hat{C}_t \odot i_t && \text{ (a modulated additive update for context)} \\ h_t &= \tanh(C_t) \odot o_t && \text{ (transform context into state and selectively output)}\end{aligned}$$

- Note: \odot represents elementwise vector product. Also, state updates now additive, not multiplicative. Training using backpropagation through time.
- Many variants of LSTM exists, e.g., using C_{t-1} in local computations, Gated Recurrent Units (GRU), etc. Mostly minor variations of basic LSTM above

Neural Nets for Unsupervised Learning

Autoencoder

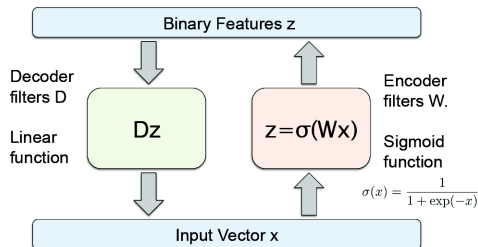
- A neural net for unsupervised feature extraction
- Basic principle: Learns an encoding of the inputs so as to recover the original input from the encodings as well as possible



- Also used to initialize deep learning models (layer-by-layer pre-training)

Autoencoder: An Example

- Real-valued inputs, binary-valued encodings



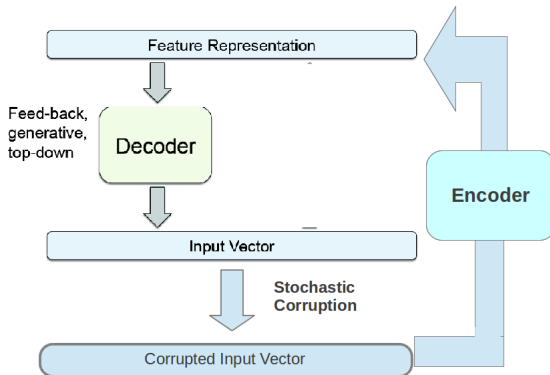
- Sigmoid encoder (parameter matrix W), linear decoder (parameter matrix D), learned via:

$$\arg \min_{D, W} E(D, W) = \sum_{n=1}^N \|D \mathbf{z}_n - \mathbf{x}_n\|^2 = \sum_{n=1}^N \|D \sigma(W \mathbf{x}_n) - \mathbf{x}_n\|^2$$

- If encoder is also linear, then autoencoder is equivalent to PCA

Denoising Autoencoders

- Idea: introduce stochastic corruption to the input; e.g.:
 - Hide some features
 - Add gaussian noise



Summary

- Looked at feedforward neural networks and extensions such as CNN
- Looked at (deep) neural nets (RNN/LSTM) for learning from sequential data
 - Methods like RNN and LSTM are widely used for learning from such data
 - Modeling and retaining context is important when modeling sequential data (desirable to have a “memory module” of some sort as in LSTMs)
- Looked at Autoencoder - Neural network for unsupervised feature extraction
- Didn't discuss some other popular methods, e.g., deep generative models, but these are based on similar underlying principles