



# Recurrent Neural Networks

---

## Deep Learning Decal

Hosted by Machine Learning at Berkeley

## Agenda

Background

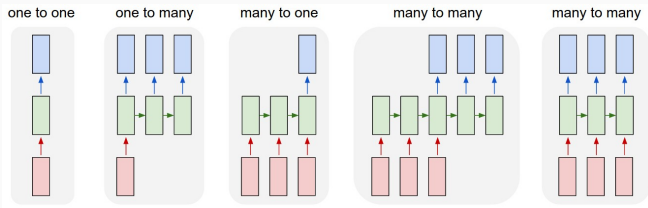
LSTM

Questions

# Background

---

- Up until now, we've only dealt with neural networks that take in a fixed input and produce a fixed output
- Recurrent neural nets (RNN) are exciting because they operate over *sequences* of vectors, which can be arbitrarily long!

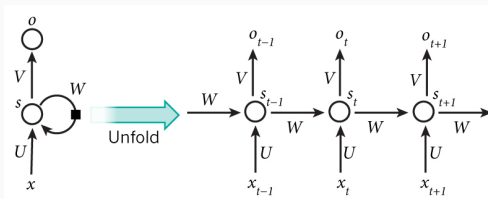


- Recurrent neural nets naturally handle tasks that involve sequences of data
  - This includes speech recognition, language modeling, translation, and image captioning
  - See “The Unreasonable Effectiveness of Recurrent Neural Networks” by Andrej Karpathy
- In practice, two slightly different architectures called the LSTM and GRU are usually used

# How do RNNs work?



- $x_t$  is the input at time step  $t$
- $s_t$  is the hidden state at time step  $t$
- $o_t$  is the hidden state at time step  $t$



- The hidden states capture information from earlier inputs, it serves as the “memory” of the network.
  - $s_t = f(Ux_t + Ws_{t-1})$  Common choices for  $f$  include tanh or ReLU.
- $o_t$  is usually some function of the hidden state
  - $\text{softmax}(Vs_t)$  generates a vector of probabilities

- You can think of recurrent neural networks as many copies of the same network, each passing a message to its successor.
  - This reduces the number of parameters we need to learn
- We don't necessarily need inputs and outputs at every time step—the key feature that RNNs share is the hidden state.

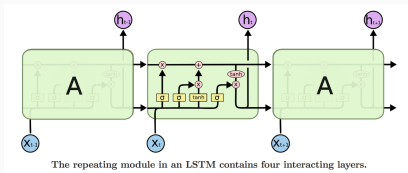


- We can use the same optimization techniques as before!
  - Define a loss function and use back propagation to update weights
- In practice, vanishing gradients can be an issue in training RNNs
  - idea: Updates to earlier layers are multiplied by many numbers  $< 1$  so they vanish
  - Using LSTMs or GRUs is one way around this problem

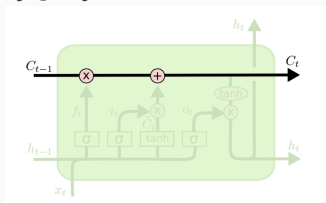
# LSTM

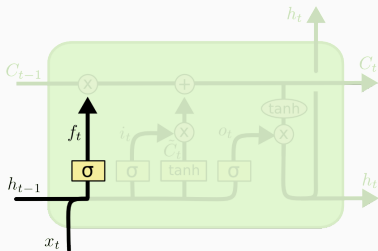
---

- Vanishing gradient problem
- Handling long-term dependencies
  - Example: In language modeling, the neural net often needs to remember information from several words ago. "I grew up in France ... I speak fluent (blank)"

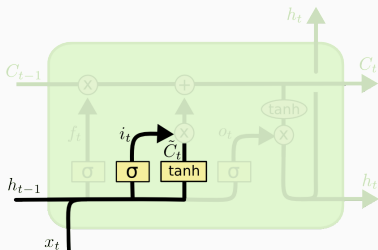


- Cell state is like a conveyor belt, gets updated via linear interactions
- It's for information to flow unchanged; information is added via regulated gates
- $C_t = f_t \odot C_{t-1} + i_t \odot \tilde{C}_t$



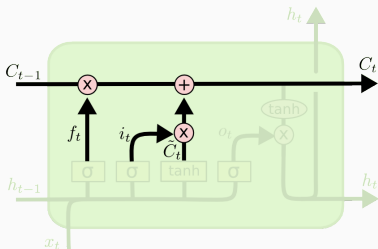


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

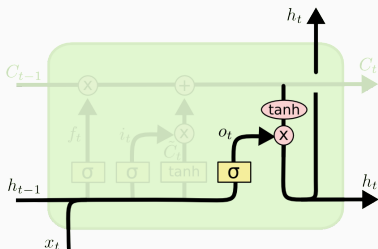


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$



# What do cell states do?



Cell sensitive to position in line:

The sole importance of the crossing of the Berezina lies in the fact that it plainly and indubitably proved the fallacy of all the plans for cutting off the enemy's retreat and the soundness of the only possible line of action--the one Kutuzov and the general mass of the army demanded--namely, simply to follow the enemy up. The French crowd fled at a continually increasing speed and all its energy was directed to reaching its goal. It fled like a wounded animal and it was impossible to block its path. This was shown not so much by the arrangements it made for crossing as by what took place at the bridges. When the bridges broke down, unarmed soldiers, people from Moscow and women with children who were with the French transport, all--carried on by vis inertiae--pressed forward into boats and into the ice-covered water and did not, surrender.

Cell that turns on inside quotes:

"You mean to imply that I have nothing to eat out of.... On the contrary, I can supply you with everything even if you want to give dinner parties," warmly replied Chichagov, who tried by every word he spoke to prove his own rectitude and therefore imagined Kutuzov to be animated by the same desire.

Kutuzov, shrugging his shoulders, replied with his subtle penetrating smile: "I meant merely to say what I said."

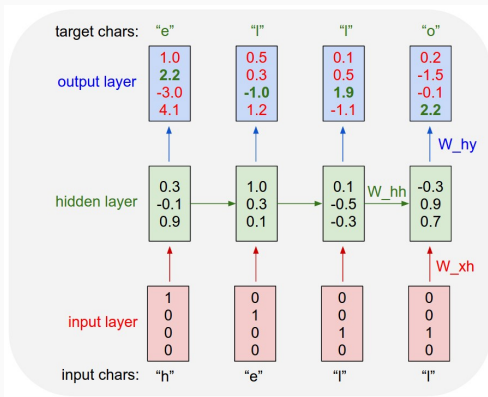
Cell that robustly activates inside if statements:

```
static int __dequeue_signal(struct sigpending *pending, sigset_t *mask,
                           siginfo_t *info)
{
    int sig = next_signal(pending, mask);
    if (sig) {
        if (current->notifier) {
            if (sigismember(current->notifier_mask, sig)) {
                if (!(current->notifier)(current->notifier_data)) {
                    clear_thread_flag(TIF_SIGPENDING);
                    return 0;
                }
            }
        }
        collect_signal(sig, pending, info);
    }
    return sig;
}
```

A large portion of cells are not easily interpretable. Here is a typical example:

```
/* Unpack a filter field's string representation from user-space
 * buffer. */
char *audit_unpack_string(void **bufp, size_t *remain, size_t len)
{
    if
```

- Many different LSTM variants, they all work about the same
  - See *LSTM: A Search Space Odyssey*, Greff et al. for more!



- Many publicly available implementations, mine is at <https://github.com/michaelrzhang/Char-RNN>

For  $\bigoplus_{i=1,\dots,m} \mathcal{L}_{i,\bullet}$  where  $\mathcal{L}_{i,\bullet} = 0$ , hence we can find a closed subset  $\mathcal{H}$  in  $\mathcal{H}$  and any sets  $\mathcal{F}$  on  $X$ ,  $U$  is a closed immersion of  $S$ , then  $U \rightarrow T$  is a separated algebraic space.

*Proof.* Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparably in the fibre product covering we have to prove the lemma generated by  $\prod Z \times_U U \rightarrow V$ . Consider the maps  $M$  along the set of points  $\text{Sch}/_{\text{fppf}}$  and  $U \rightarrow U$  is the fibre category of  $S$  in  $U$  in Section, ?? and the fact that any  $U$  affine, see Morphisms, Lemma ?? Hence we obtain a scheme  $S$  and any open subset  $W \subset U$  in  $\text{Sh}(G)$  such that  $\text{Spec}(R) \rightarrow S$  is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that  $f_i$  is of finite presentation over  $S$ . We claim that  $\mathcal{O}_{X,x}$  is a scheme where  $x, x', x'' \in S'$  such that  $\mathcal{O}_{X,x'} \rightarrow \mathcal{O}'_{X',x'}$  is separated. By Algebra, Lemma ?? we can define a map of complexes  $\text{GL}_{\mathcal{O}}(x'/S')$  and we win.  $\square$

To prove study we see that  $\mathcal{F}_i^0$  is a covering of  $X'$ , and  $\mathcal{T}_i$  is an object of  $\mathcal{F}_{X/S}$  for  $i > 0$  and  $\mathcal{F}_i$  exists and let  $\mathcal{F}_i$  be a product of  $\mathcal{O}_{X/S}$ -modules on  $\mathcal{C}$  as a  $\mathcal{F}$ -module. In particular  $\mathcal{F} = U/\mathcal{F}$  we have to show that

$$\tilde{\mathcal{H}}^* = \mathcal{T}^* \otimes_{\text{Hom}(S)} \mathcal{O}_{S,S} - \mathcal{I}_X^{-1}(\mathcal{F})$$

is a unique morphism of algebraic stacks. Note that

$$\text{Arrows} = (\text{Sch}/S)_{\text{fppf}}^{\text{op}}/(\text{Sch}/S)_{\text{fppf}}$$

and

$$V = \Gamma(S, \mathcal{O}) \rightarrow (U, \text{Spec}(A))$$

is an open subset of  $X$ . Thus  $U$  is affine. This is a continuous map of  $X$  is the inverse, the gropsoid scheme  $S$ .

*Proof.* See discussion of sheaves of sets.  $\square$

The result for prove any open covering follows from the loss of Example ?? It may replace  $S$  by  $X_{\text{pro-étale}}$  which gives an open subspace of  $X$  and  $T$  equal to  $S_{\text{pro-étale}}$  see Descent, Lemma ?? Namely, by Lemma ?? we see that  $R$  is geometrically regular over  $S$ .

**Lemma 0.1.** Assume (3) and (3) by the construction in the description.

Suppose  $X = \lim[X]$  (by the formal open covering  $X$  and a single map  $\text{Proj}_X(\mathcal{A}) = \text{Spec}(B)$  over  $U$  compatible with the complex

$$\text{Set}(\mathcal{A}) = \Gamma(X, \mathcal{O}_{X, \mathcal{O}_X})$$

When in this case of to show that  $\mathcal{Q} \rightarrow \mathcal{C}_{Z/X}$  is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are arbitrary. If  $T$  is surjective we may assume that  $T$  is connected with residue fields of  $S$ . Moreover there exists a closed subspace  $Z \subset X$  of  $X$  where  $U$  in  $X'$  is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem

(1)  $f$  is locally of finite type. Since  $S = \text{Spec}(R)$  and  $Y = \text{Spec}(R)$ .

*Proof.* This is form all sheaves of sheaves on  $X$ . But given a scheme  $U$  and a surjective étale morphism  $U \rightarrow X$ . Let  $U \cap U = \prod_{i=1,\dots,r} U_i$  be the scheme  $X$  over  $S$  at the schemes  $X_i \rightarrow X$  and  $U = \lim_i X_i$ .  $\square$

The following lemma surjective restrocomposes of this implies that  $\mathcal{F}_{\pi_0} = \mathcal{F}_{\pi_0} = \mathcal{F}_{X \dots \beta}$ .

**Lemma 0.2.** Let  $X$  be a locally Noetherian scheme over  $S$ ,  $E = \mathcal{F}_{X/S}$ . Set  $\mathcal{I} = \mathcal{I}_i \subset \mathcal{I}_n$ . Since  $\mathcal{I}^n \subset \mathcal{I}^n$  are nonzero over  $i_0 \leq p$  is a subset of  $\mathcal{J}_{n,0} \circ \mathcal{A}_2$  works.

**Lemma 0.3.** In Situation ?? Hence we may assume  $\mathfrak{q}' = 0$ .

*Proof.* We will use the property we see that  $p$  is the next functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_X) = \mathcal{O}_X(D)$$

where  $K$  is an  $F$ -algebra where  $\delta_{n+1}$  is a scheme over  $S$ .  $\square$

- Works well with many different data sets!

## Questions

---

Questions?