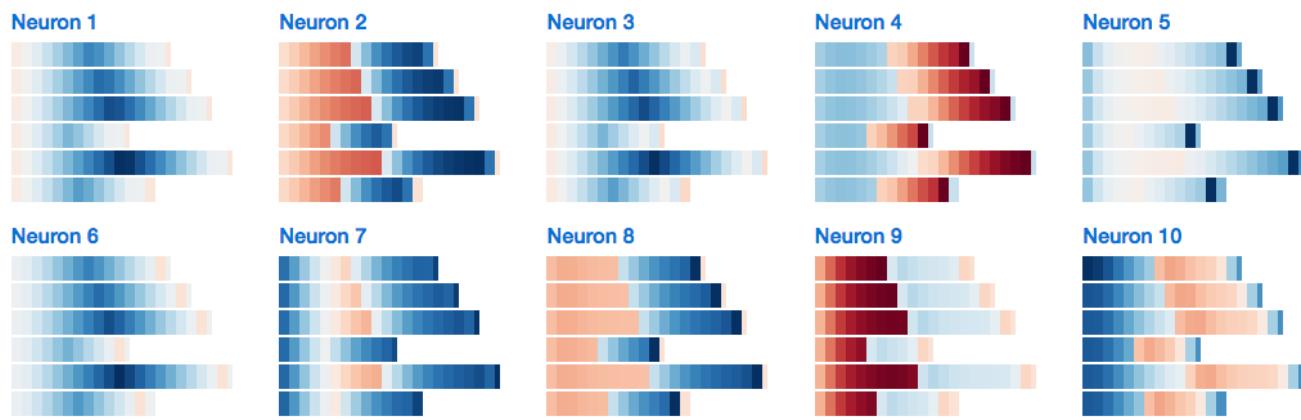


# Exploring LSTMs

LSTMs are behind a lot of the amazing achievements deep learning has made in the past few years, and they're a fairly simple extension to neural networks under the right view. So I'll try to present them as intuitively as possible – in such a way that you could have discovered them yourself.

But first, a picture:



Aren't LSTMs beautiful? Let's go.

(Note: if you're already familiar with neural networks and LSTMs, skip to the middle – the first half of this post is a tutorial.)

## Neural Networks

Imagine we have a sequence of images from a movie, and we want to label each image with an activity (is this a fight?, are the characters talking?, are the characters eating?).

How do we do this?

One way is to ignore the sequential nature of the images, and build a per-image classifier that considers each image in isolation. For example, given enough images and labels:

- Our algorithm might first learn to detect low-level patterns like shapes and edges.
- With more data, it might learn to combine these patterns into more complex ones, like faces (two circular things atop a triangular thing atop an oval thing) or cats.
- And with even more data, it might learn to map these higher-level patterns into activities themselves (scenes with mouths, steaks, and forks are probably about eating).

This, then, is a deep neural network: it takes an image input, returns an activity output, and – just as we might learn to detect patterns in puppy behavior without knowing

anything about dogs (after seeing enough corgis, we discover common characteristics like fluffy butts and drumstick legs; next, we learn advanced features like splotting) – in between it learns to represent images through hidden layers of representations.

## Mathematically

I assume people are familiar with basic neural networks already, but let's quickly review them.

- A neural network with a single hidden layer takes as input a vector  $x$ , which we can think of as a set of neurons.
- Each input neuron is connected to a hidden layer of neurons via a set of learned weights.
- The  $j$ th hidden neuron outputs  $h_j = \phi(\sum_i w_{ij}x_i)$ , where  $\phi$  is an activation function.
- The hidden layer is fully connected to an output layer, and the  $j$ th output neuron outputs  $y_j = \sum_i v_{ij}h_i$ . If we need probabilities, we can transform the output layer via a **softmax** function.

In matrix notation:

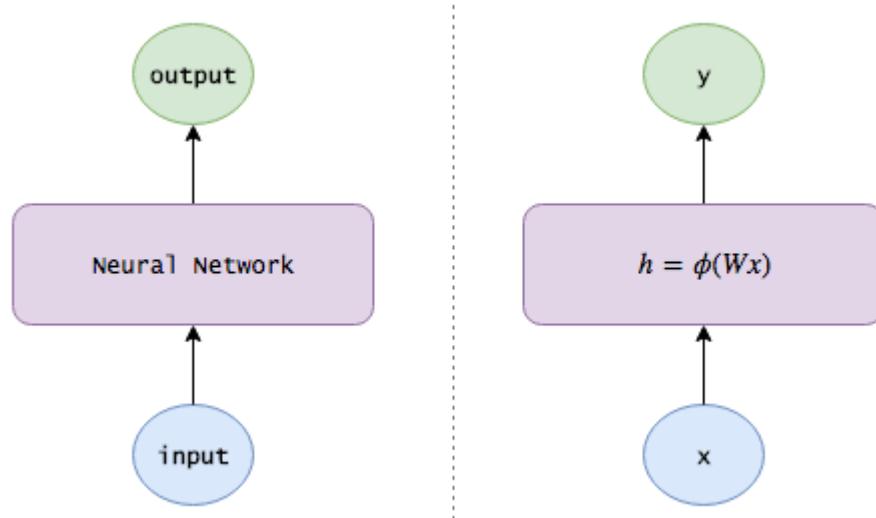
$$h = \phi(Wx)$$

$$y = Vh$$

where

- $x$  is our input vector
- $W$  is a weight matrix connecting the input and hidden layers
- $V$  is a weight matrix connecting the hidden and output layers
- Common activation functions for  $\phi$  are the **sigmoid function**,  $\sigma(x)$ , which squashes numbers into the range  $(0, 1)$ ; the **hyperbolic tangent**,  $\tanh(x)$ , which squashes numbers into the range  $(-1, 1)$ , and the **rectified linear unit**,  $ReLU(x) = \max(0, x)$ .

Here's a pictorial view:



(Note: to make the notation a little cleaner, I assume  $x$  and  $h$  each contain an extra bias neuron fixed at 1 for learning bias weights.)

# Remembering Information with RNNs

Ignoring the sequential aspect of the movie images is pretty ML 101, though. If we see a scene of a beach, we should boost beach activities in future frames: an image of someone in the water should probably be labeled swimming, not bathing, and an image of someone lying with their eyes closed is probably suntanning. If we remember that Bob just arrived at a supermarket, then even without any distinctive supermarket features, an image of Bob holding a slab of bacon should probably be categorized as shopping instead of cooking.

So what we'd like is to let our model track the state of the world:

1. After seeing each image, the model outputs a label and also updates the knowledge it's been learning. For example, the model might learn to automatically discover and track information like location (are scenes currently in a house or beach?), time of day (if a scene contains an image of the moon, the model should remember that it's nighttime), and within-movie progress (is this image the first frame or the 100th?). Importantly, just as a neural network automatically discovers hidden patterns like edges, shapes, and faces without being fed them, our model should automatically discover useful information by itself.
2. When given a new image, the model should incorporate the knowledge it's gathered to do a better job.

This, then, is a recurrent neural network. Instead of simply taking an image and returning an activity, an RNN also maintains internal memories about the world (weights assigned to different pieces of information) to help perform its classifications.

## Mathematically

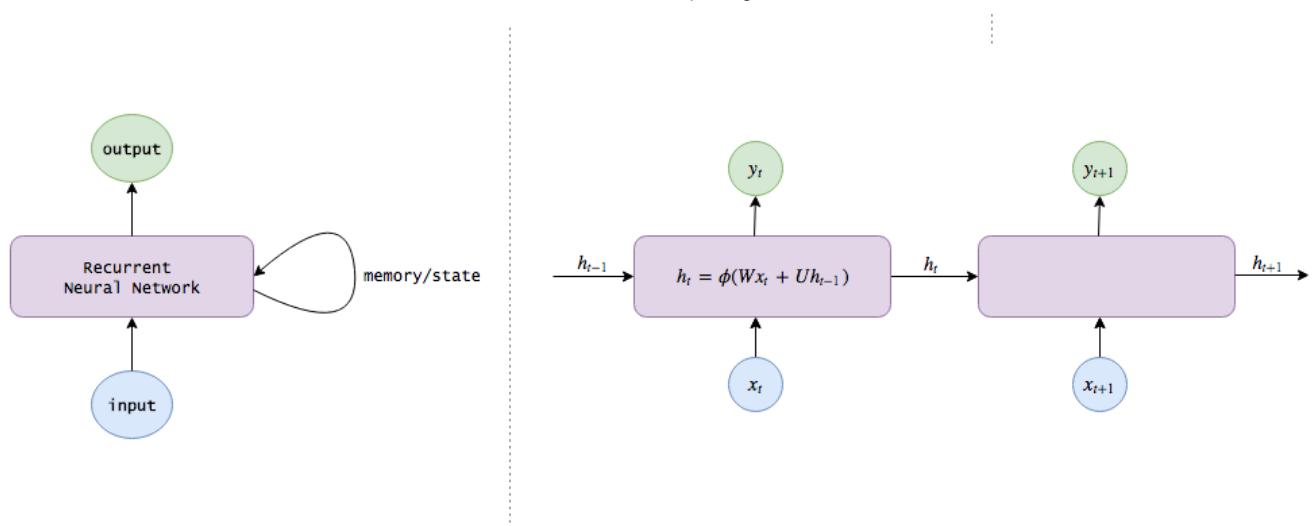
So let's add the notion of internal knowledge to our equations, which we can think of as pieces of information that the network maintains over time.

But this is easy: we know that the hidden layers of neural networks already encode useful information about their inputs, so why not use these layers as the memory passed from one time step to the next? This gives us our RNN equations:

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

$$y_t = Vh_t$$

Note that the hidden state computed at time  $t$  ( $h_t$ , our internal knowledge) is fed back at the next time step. (Also, I'll use concepts like hidden state, knowledge, memories, and beliefs to describe  $h_t$  interchangeably.)



## Longer Memories through LSTMs

Let's think about how our model updates its knowledge of the world. So far, we've placed no constraints on this update, so its knowledge can change pretty chaotically: at one frame it thinks the characters are in the US, at the next frame it sees the characters eating sushi and thinks they're in Japan, and at the next frame it sees polar bears and thinks they're on Hydra Island. Or perhaps it has a wealth of information to suggest that Alice is an investment analyst, but decides she's a professional assassin after seeing her cook.

This chaos means information quickly transforms and vanishes, and it's difficult for the model to keep a long-term memory. So what we'd like is for the network to learn how to update its beliefs (scenes without Bob shouldn't change Bob-related information, scenes with Alice should focus on gathering details about her), in a way that its knowledge of the world evolves more gently.

This is how we do it.

1. **Adding a forgetting mechanism.** If a scene ends, for example, the model should forget the current scene location, the time of day, and reset any scene-specific information; however, if a character dies in the scene, it should continue remembering that he's no longer alive. Thus, we want the model to learn a separate forgetting/remembering mechanism: when new inputs come in, it needs to know which beliefs to keep or throw away.
2. **Adding a saving mechanism.** When the model sees a new image, it needs to learn whether any information about the image is worth using and saving. Maybe your mom sent you an article about the Kardashians, but who cares?
3. So when a new input comes in, the model first forgets any long-term information it decides it no longer needs. Then it learns which parts of the new input are worth using, and saves them into its long-term memory.
4. **Focusing long-term memory into working memory.** Finally, the model needs to learn which parts of its long-term memory are immediately useful. For example, Bob's age may be a useful piece of information to keep in the long term (children are more likely to be crawling, adults are more likely to be working), but is probably irrelevant if he's not in the current

scene. So instead of using the full long-term memory all the time, it learns which parts to focus on instead.

This, then, is an long short-term memory network. Whereas an RNN can overwrite its memory at each time step in a fairly uncontrolled fashion, an LSTM transforms its memory in a very precise way: by using specific learning mechanisms for which pieces of information to remember, which to update, and which to pay attention to. This helps it keep track of information over longer periods of time.

## Mathematically

Let's describe the LSTM additions mathematically.

At time  $t$ , we receive a new input  $x_t$ . We also have our long-term and working memories passed on from the previous time step,  $ltm_{t-1}$  and  $wm_{t-1}$  (both n-length vectors), which we want to update.

We'll start with our long-term memory. First, we need to know which pieces of long-term memory to continue remembering and which to discard, so we want to use the new input and our working memory to learn a **remember gate** of n numbers between 0 and 1, each of which determines how much of a long-term memory element to keep. (A 1 means to keep it, a 0 means to forget it entirely.)

Naturally, we can use a small neural network to learn this remember gate:

$$\text{remember}_t = \sigma(W_r x_t + U_r w_{m,t-1})$$

(Notice the similarity to our previous network equations; this is just a shallow neural network. Also, we use a sigmoid activation because we need numbers between 0 and 1.)

Next, we need to compute the information we can learn from  $x_t$ , i.e., a **candidate addition to our long-term memory**:

$$ltm'_t = \phi(W_l x_t + U_l w_{m,t-1})$$

$\phi$  is an activation function, commonly chosen to be  $\tanh$ .

Before we add the candidate into our memory, though, we want to learn which parts of it are actually worth using and saving:

$$\text{save}_t = \sigma(W_s x_t + U_s w_{m,t-1})$$

(Think of what happens when you read something on the web. While a news article might contain information about Hillary, you should ignore it if the source is Breitbart.)

Let's now combine all these steps. After forgetting memories we don't think we'll ever need again and saving useful pieces of incoming information, we have our updated long-term memory:

$$ltm_t = remember_t \circ ltm_{t-1} + save_t \circ ltm'_t$$

where  $\circ$  denotes element-wise multiplication.

Next, let's update our working memory. We want to learn how to focus our long-term memory into information that will be immediately useful. (Put differently, we want to learn what to move from an external hard drive onto our working laptop.) So we learn a focus/attention vector:

$$focus_t = \sigma(W_f x_t + U_f w m_{t-1})$$

Our working memory is then

$$w m_t = focus_t \circ \phi(ltm_t)$$

In other words, we pay full attention to elements where the focus is 1, and ignore elements where the focus is 0.

And we're done! Hopefully this made it into your long-term memory as well.

To summarize, whereas a vanilla RNN uses one equation to update its hidden state/memory:

$$h_t = \phi(W x_t + U h_{t-1})$$

An LSTM uses several:

$$ltm_t = remember_t \circ ltm_{t-1} + save_t \circ ltm'_t$$

$$w m_t = focus_t \circ \tanh(ltm_t)$$

where each memory/attention sub-mechanism is just a mini brain of its own:

$$remember_t = \sigma(W_r x_t + U_r w m_{t-1})$$

$$save_t = \sigma(W_s x_t + U_s w m_{t-1})$$

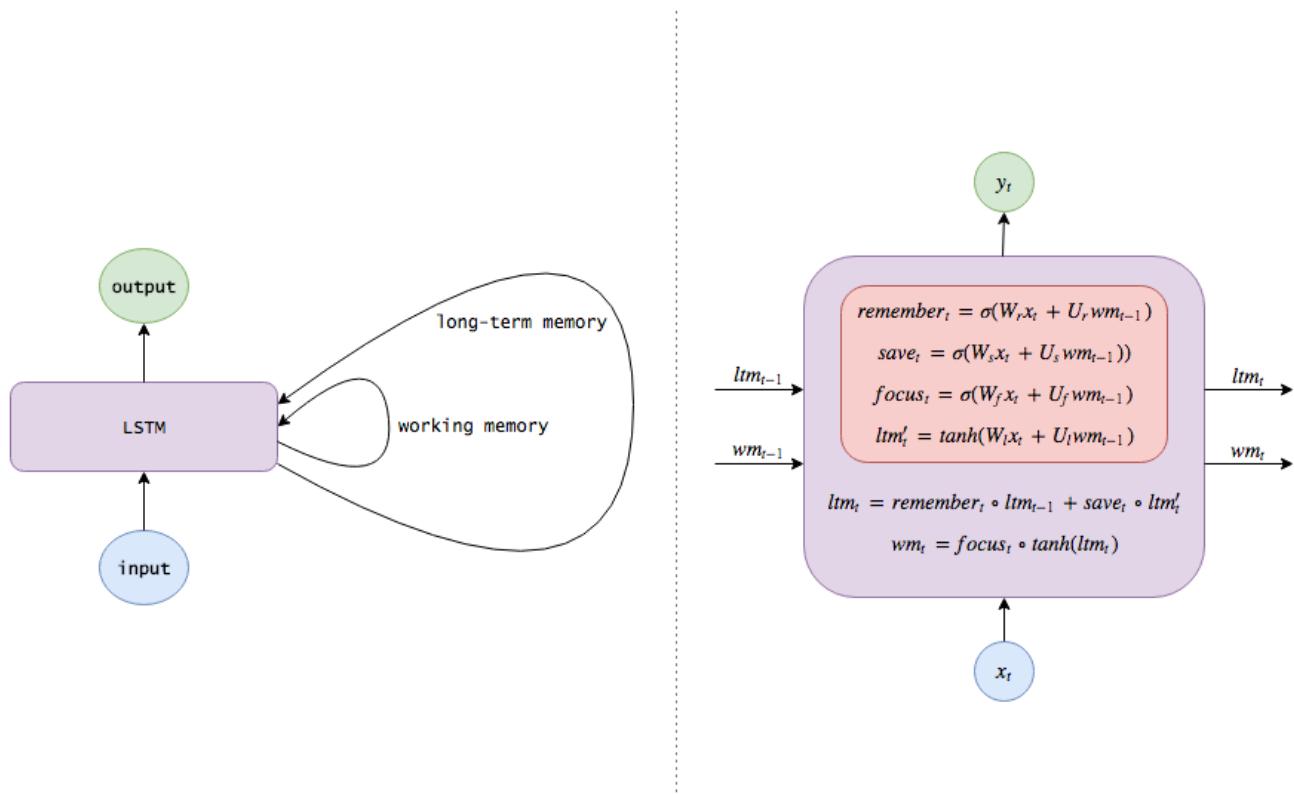
$$focus_t = \sigma(W_f x_t + U_f w m_{t-1})$$

$$ltm'_t = \tanh(W_l x_t + U_l w m_{t-1})$$

(Note: the terminology and variable names I've been using are different from the usual literature. Here are the standard names, which I'll use interchangeably from now on:

- The long-term memory,  $ltm_t$ , is usually called the **cell state**, denoted  $c_t$ .
- The working memory,  $w m_t$ , is usually called the **hidden state**, denoted  $h_t$ . This is analogous to the hidden state in vanilla RNNs.
- The remember vector,  $remember_t$ , is usually called the **forget gate** (despite the fact that a 1 in the forget gate still means to keep the memory and a 0 still means to forget it), denoted  $f_t$ .

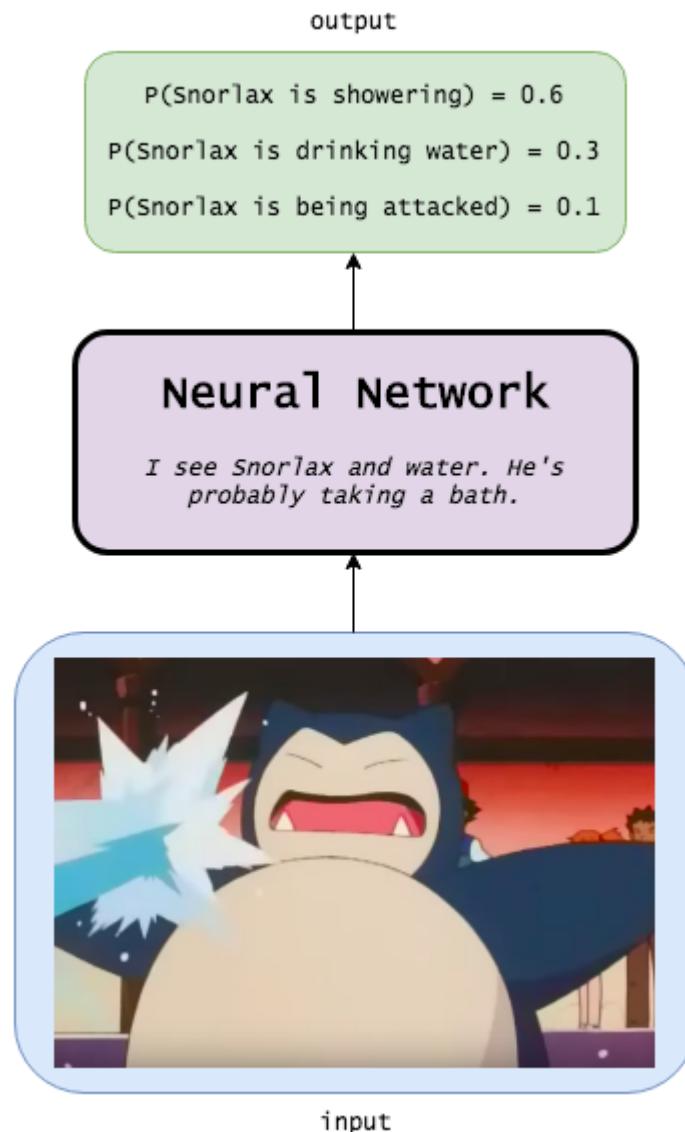
- The save vector,  $save_t$ , is usually called the input gate (as it determines how much of the input to let into the cell state), denoted  $i_t$ .
- The focus vector,  $focus_t$ , is usually called the output gate, denoted  $o_t$ . )



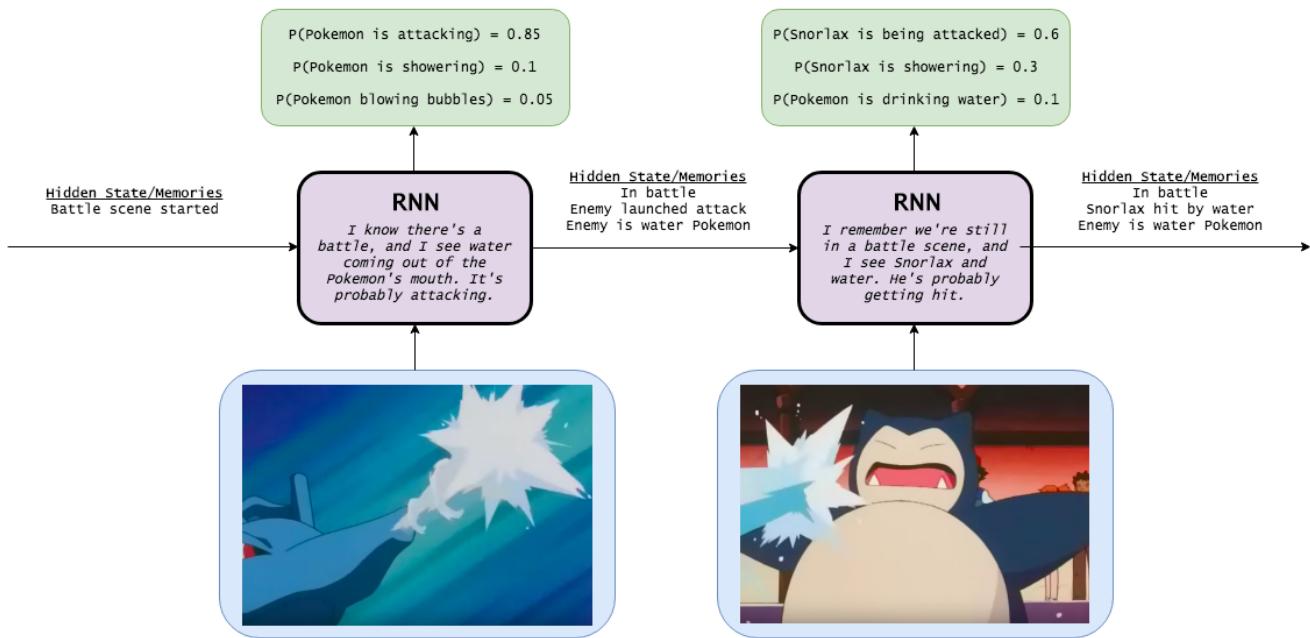
## Snorlax

I could have caught a hundred Pidgeys in the time it took me to write this post, so here's a cartoon.

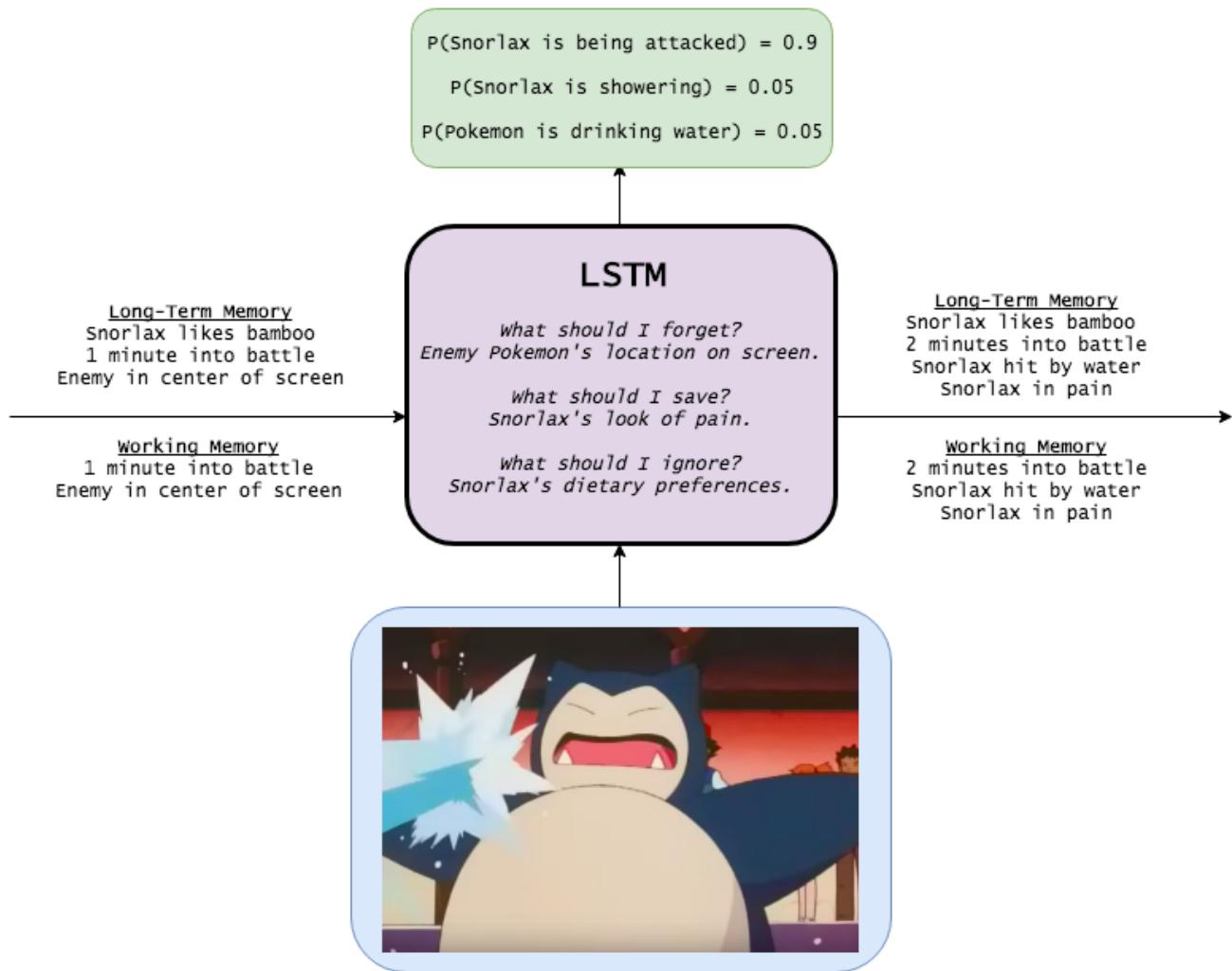
## Neural Networks



## Recurrent Neural Networks



# LSTMs



## Learning to Code

Let's look at a few examples of what an LSTM can do. Following Andrej Karpathy's [terrific post](#), I'll use character-level LSTM models that are fed sequences of characters and trained to predict the next character in the sequence.

While this may seem a bit toyish, character-level models can actually be very useful, even on top of word models. For example:

- Imagine a code autocomplete smart enough to allow you to program on your phone. An LSTM could (in theory) track the return type of the method you're currently in, and better suggest which variable to return; it could also know without compiling whether you've made a bug by returning the wrong type.
- NLP applications like machine translation often have trouble dealing with rare terms. How do you translate a word you've never seen before, or convert adjectives to adverbs? Even if you know what a tweet means, how do you generate a new hashtag to capture it? Character models can daydream new terms, so this is another area with [interesting applications](#).

So to start, I spun up an EC2 p2.xlarge spot instance, and trained a 3-layer LSTM on the [Apache Commons Lang codebase](#). Here's a program it generates after a few hours.

```
1  /*
2   * Licensed to the Apache Software Foundation (ASF) under one or more
3   * contributor license agreements. See the NOTICE file distributed with
4   * this work for additional information regarding copyright ownership.
5   * The ASF licenses this file to You under the Apache License, Version 2.0
6   * (the "License"); you may not use this file except in compliance with
7   * the License. You may obtain a copy of the License at
8   *
9   *      http://www.apache.org/licenses/LICENSE-2.0
10  *
11  * Unless required by applicable law or agreed to in writing, software
12  * distributed under the License is distributed on an "AS IS" BASIS,
13  * WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
14  * See the License for the specific language governing permissions and
15  * limitations under the License.
16  */
17
18 package org.apache.commons.math4.linear;
19
20 import java.text.NumberFormat;
21 import java.io.ByteArrayInputStream;
22 import java.io.ObjectOutputStream;
23 import java.io.ObjectInputStream;
24 import java.util.ArrayList;
25 import java.util.List;
26
27 import org.apache.commons.math4.optim.nonlinear.scalar.GoalType;
28 import org.apache.commons.math4.ml.neuralnet.sofm.NeuronSquareMesh2D;
29 import org.apache.commons.math4.distribution.DescriptiveStatistics;
30 import org.apache.commons.math4.optim.nonlinear.scalar.NodeFieldIntegrator;
31 import org.apache.commons.math4.optim.nonlinear.scalar.GradientFunction;
32 import org.apache.commons.math4.optim.PointValuePair;
33 import org.apache.commons.numbers.core.Precision;
34
35 /**
36  * <p>Natural infinite is defined in basic eigenvalues of a transform are in a subconsider fo
37  *
38  * <p>This implementation is the computation at a collection of a set of the solvers.</p>
39  * <p>
40  * This class is returned the default precision parameters after a new value for the interpol
41  * <p>
```

```
42 * The distribution values do not ratio example function containing this interface, which sh
43 * <p>
44 * This class generates a new standard deviation of the following conventions, the variance w
45 * constructor, and invoke the interpolation arrays</li>
46 * <li>{@code a < 1} and {@code this} the regressions returned by calling
47 * the same special corresponding to a representation.
48 * </p>
49 *
50 * @since 1.2
51 */
52 public class SinoutionIntegrator implements Serializable {
53
54     /** Serializable version identifier */
55     private static final long serialVersionUID = -7989543519820244888L;
56
57     /**
58      * Start distance between the instance and a result (does not all lead to the number of s
59      * <p>
60      * Note that this implementation this can prevent the permutation of the preneved statist
61      * </p>
62      * <p>
63      * <strong>Preconditions</strong>: <ul>
64      * <li>Returns number of samples and the designated subarray, or
65      * if it is null, {@code null}. It does not dofine the base number.</p>
66      *
67      * @param source the number of left size of the specified value
68      * @param numberPoints number of points to be checked
69      * @return the parameters for a public function.
70     */
71     public static double fitness(final double[] sample) {
72         double additionalComputed = Double.POSITIVE_INFINITY;
73         for (int i = 1; i < dim; i++) {
74             final double coefficients[i] = point[i] * coefficients[i];
75             double diff = a * FastMath.cos(point[i]);
76             final double sum = FastMath.max(random.nextDouble(), alpha);
77             final double sum = FastMath.sin(optimal[i].getReal() - cholenghat);
78             final double lower = gamma * cHessian;
79             final double fs = factor * maxIterationCount;
80             if (temp > numberPoints - 1) {
81                 final int pma = points.size();
82                 boolean partial = points.toString();
83                 final double segments = new double[2];
84                 final double sign = pti * x2;
85                 double n = 0;
```

```
86         for (int i = 0; i < n; i++) {
87             final double ds = normalizedState(i, k, difference * factor);
88             final double inv = alpha + temp;
89             final double rsigx = FastMath.sqrt(max);
90             return new String(degree, e);
91         }
92     }
93     // Perform the number to the function parameters from one count of the values
94     final PointValuePair part = new PointValuePair[n];
95     for (int i = 0; i < n; i++) {
96         if (i == 1) {
97             numberPoints = 1;
98         }
99         final double dev = FastMath.log(perturb(g, norm), values[i]);
100        if (Double.isNaN(y) &&
101            NaN) {
102            sum /= samples.length;
103        }
104        double i = 1;
105        for (int i = 0; i < n; i++) {
106            statistics[i] = FastMath.abs(point[i].sign() + rhs[i]);
107        }
108        return new PointValuePair(true, params);
109    }
110}
111}
112
113 /**
114 * Computes the number of values
115 * @throws NotPositiveException if {@code NumberIsTooSmallException} if {@code seed <= 0}
116 * @throws NullArgumentException if row or successes is null
117 */
118 public static double numericalMean(double value) {
119     if (variance == null) {
120         throw new NotStrictlyPositiveException(LocalizedFormats.NUMBER_OF_SUBCORSE_TRANSI-
121                                         p, numberSuccesses, true);
122     }
123     return sum;
124 }
125
126 /**
127 * {@inheritDoc}
128 */
129 @Override
```

```

130     public LeastSquaresProblem create(final StatisticalSummary sampleStats1,
131                                         final double[] values, final double alpha) throws Math
132         final double sum = sumLogImpl.toSubSpace(sample);
133         final double relativeAccuracy = getSumOfLogs();
134         final double[] sample1 = new double[dimension];
135
136         for (int i = 0; i < result.length; i++) {
137             verifyInterval.solve(params, alpha);
138         }
139         return max;
140     }
141
142     /**
143      * Test creates a new PolynomialFunction function
144      * @see #applyTo(double)
145      */
146     @Test
147     public void testCosise() {
148         final double p = 7.7;
149         final double expected = 0.0;
150         final SearchInterval d = new Power(1.0, 0.0);
151         final double penalty = 1e-03;
152         final double init = 0.245;
153         final double t = 0.2;
154         final double result = (x + 1.0) / 2.0;
155         final double numeratorAdd = 13;
156         final double bhigh = 2 * (k - 1) * Math.acos();
157
158         Assert.assertEquals(0.0, true);
159         Assert.assertTrue(percentile.evaluate(singletonArray), 0);
160         Assert.assertEquals(0.0, getNumberOfTrials(0, 0), 1E-10);
161         Assert.assertEquals(0.201949230731, percentile.evaluate(specialValues), 1.0e-3);
162         Assert.assertEquals(-10.0, distribution.inverseCumulativeProbability(0.50), 0);
163         Assert.assertEquals(0.0, solver.solve(100, f, 1.0, 0.5), 1.0e-10);
164     }

```

While the code certainly isn't perfect, it's better than a lot of data scientists I know. And we can see that the LSTM has learned a lot of interesting (and correct!) coding behavior:

- It knows how to structure classes: a license up top, followed by packages and imports, followed by comments and a class definition, followed by variables and methods. Similarly, it knows how to create methods: comments follow the correct orders (description, then @param, then @return, etc.), decorators are properly placed, and non-void methods end with

appropriate return statements. Crucially, this behavior spans long ranges of code – see how giant the blocks are!

- It can also track subroutines and nesting levels: indentation is always correct, and if statements and for loops are always closed out.
- It even knows how to create tests.

How does the model do this? Let's look at a few of the hidden states.

Here's a neuron that seems to track the code's outer level of indentation:

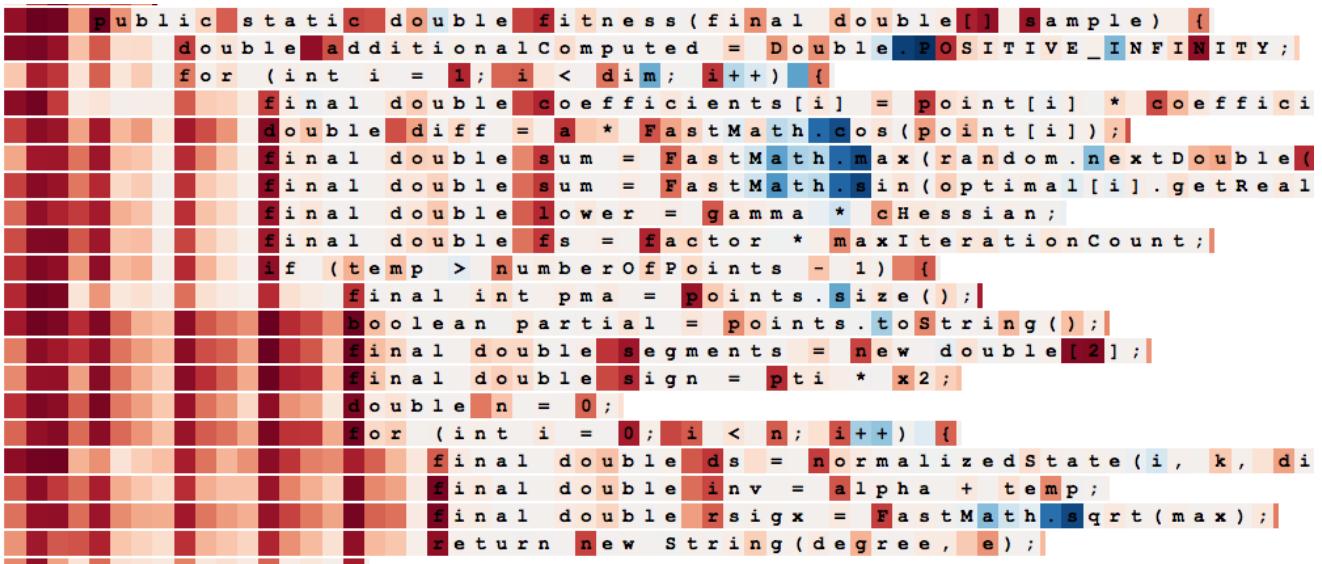
(As the LSTM moves through the sequence, its neurons fire at varying intensities. The picture represents one particular neuron, where each row is a sequence and characters are color-coded according to the neuron's intensity; dark blue shades indicate large, positive activations, and dark red shades indicate very negative activations.)

```

public static double fitness(final double[] sample) {
    double additionalComputed = Double.POSITIVE_INFINITY;
    for (int i = 1; i < dim; i++) {
        final double coefficients[i] = point[i] * coefficients[i];
        double diff = a * FastMath.cos(point[i]);
        final double sum = FastMath.max(random.nextDouble(), alpha);
        final double lower = gamma * cHessian;
        final double fs = factor * maxIterationCount;
        if (temp > numberOFPoints - 1) {
            final int pma = points.size();
            boolean partial = points.toString();
            final double segments = new double[2];
            final double sign = pti * x2;
            double n = 0;
            for (int i = 0; i < n; i++) {
                final double ds = normalizedState(i, k, difference * factor);
                final double inv = alpha + temp;
                final double rsigx = FastMath.sqrt(max);
                return new String(degree, e);
            }
        }
        // Perform the number to the function parameters from one count of the val
        final PointValuePair part = new PointValuePair[n];
        for (int i = 0; i < n; i++) {
            if (i == 1) {
                numberOFPoints = 1;
            }
            final double dev = FastMath.log(perturb(g, norm), values[i]);
            if (Double.isNaN(y) && !Double.isNaN(x)) {
                sum /= samples.length;
            }
            double i = 1;
            for (int i = 0; i < n; i++) {
                statistics[i] = FastMath.abs(point[i].sign() + rhs[i]);
            }
        }
        return new PointValuePair(true, params);
    }
}

```

And here's a neuron that counts down the spaces between tabs:



```

public static double fitness(final double[] sample) {
    double additionalComputed = Double.POSITIVE_INFINITY;
    for (int i = 1; i < dim; i++) {
        final double coefficients[i] = point[i] * coefficient;
        double diff = a * FastMath.cos(point[i]);
        final double sum = FastMath.max(random.nextDouble());
        final double sum = FastMath.sin(optimal[i].getReal());
        final double lower = gamma * cHessian;
        final double fs = factor * maxIterationCount;
        if (temp > number_of_points - 1) {
            final int pma = points.size();
            boolean partial = points.toString();
            final double segments = new double[2];
            final double sign = pti * x2;
            double n = 0;
            for (int i = 0; i < n; i++) {
                final double ds = normalizedState(i, k, di);
                final double inv = alpha + temp;
                final double rsigx = FastMath.sqrt(max);
            }
            return new String(degree, e);
        }
    }
}

```

For kicks, here's the output of a different 3-layer LSTM trained on TensorFlow's codebase:

```

1 """Tests for softplus layer tests."""
2
3 from __future__ import absolute_import
4 from __future__ import division
5 from __future__ import print_function
6
7 import collections
8 import numpy as np
9
10 from tensorflow.python.platform import test
11
12
13 class InvalidAllOpCost(Experiment):
14
15     def _runTestToIndForDead(self):
16         return self._divs()
17
18     def testPad(self):
19         with ops.Graph().as_default():
20             var = sess.run(bucketized_op)
21             self.assertAllClose(
22                 list(variables.global_variables()), status.eval())
23
24     def testHttpTimenaterRoutingOptimizerSize(self):
25         with self.test_session() as sess:
26             table = lookup_ops.IdTableWithHashBuckets(
27                 keys=['id', 'z'],
28                 example_id_column='price',
29                 num_outputs=6,

```

```

30     input_columns=['dummy_range', 'feature', 'dimensions'])
31
32     with self.assertRaisesRegexp(ValueError, 'Expected dict of rank dimensions'):
33         fc.numeric_column('aaa', indices=[[0, 0], [1, 0]], dtype=dtypes.int64)
34     output = table.lookup(input_string)
35
36     # all input tensors in SparseColumn has dimensions [end_back_prob, dimension] in the form
37     with self.assertRaisesRegexp(
38         TypeError, "Shape of values must be specified during training."):
39         fc.bucketized_column(attrs, boundaries=[62, 62])

```

There are [plenty of other fun examples](#) floating around the web, so check them out if you want to see more.

## Investigating LSTM Internals

Let's dig a little deeper. We looked in the last section at examples of hidden states, but I wanted to play with LSTM cell states and their other memory mechanisms too. Do they fire when we expect, or are there surprising patterns?

### Counting

To investigate, let's start by teaching an LSTM to count. (Remember how the Java and Python LSTMs were able to generate proper indentation!) So I generated sequences of the form

aaaaaaXbbbbbb

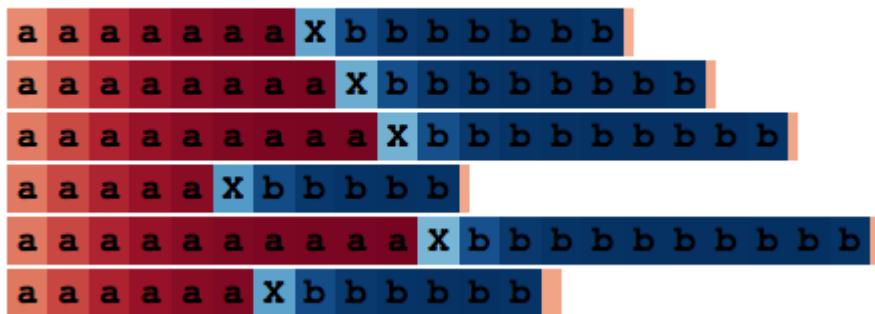
(N "a" characters, followed by a delimiter X, followed by N "b" characters, where  $1 \leq N \leq 10$ ), and trained a single-layer LSTM with 10 hidden neurons.

As expected, the LSTM learns perfectly within its training range – and can even generalize a few steps beyond it. (Although it starts to fail once we try to get it to count to 19.)

aaaaaaaaaaaaaaXbbbbbbbbbbbbbbb  
 aaaaaaaaaaaaaaaXbbbbbbbbbbbbbbb  
 aaaaaaaaaaaaaaaXbbbbbbbbbbbbbbb  
 aaaaaaaaaaaaaaaXbbbbbbbbbbbbbbb  
 aaaaaaaaaaaaaaaXbbbbbbbbbbbbbbb # Here it begins to fail: the model is give

We expect to find a hidden state neuron that counts the number of a's if we look at its internals. And we do:

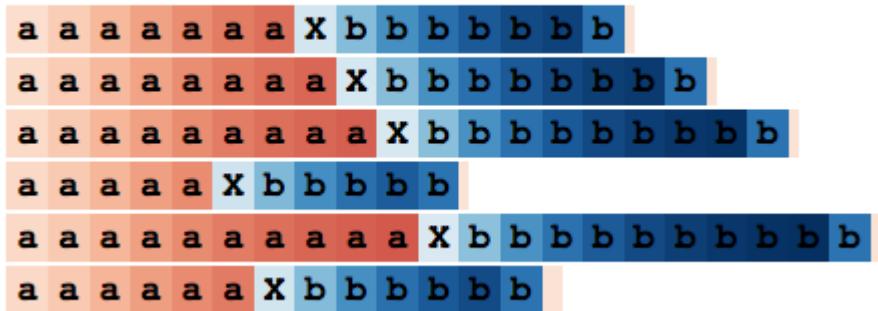
## Hidden State



I built a [small web app](#) to play around with LSTMs, and [Neuron #2](#) seems to be counting both the number of a's it's seen, as well as the number of b's. (Remember that cells are shaded according to the neuron's activation, from dark red [-1] to dark blue [+1].)

What about the cell state? It behaves similarly:

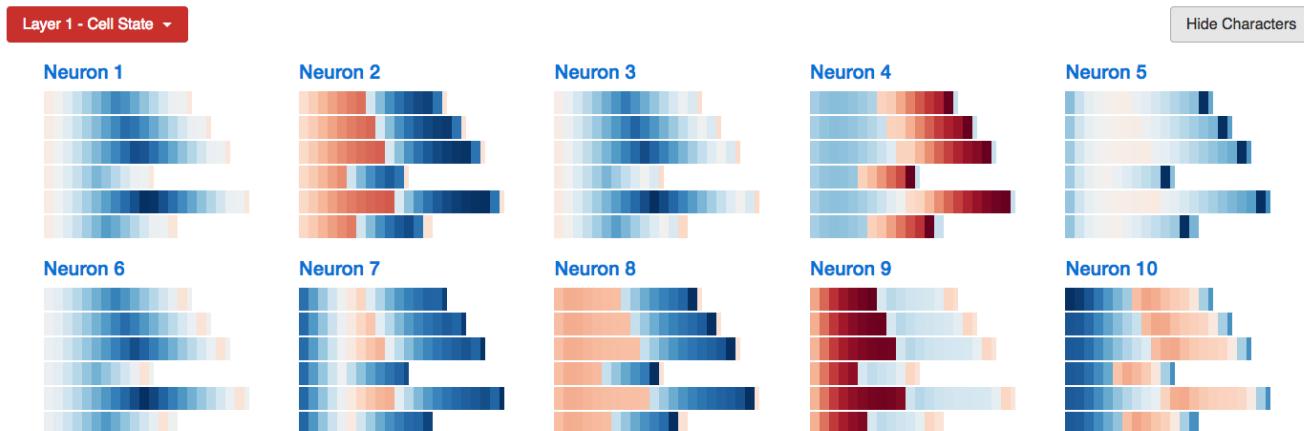
## Cell State



One interesting thing is that the working memory looks like a "sharpened" version of the long-term memory. Does this hold true in general?

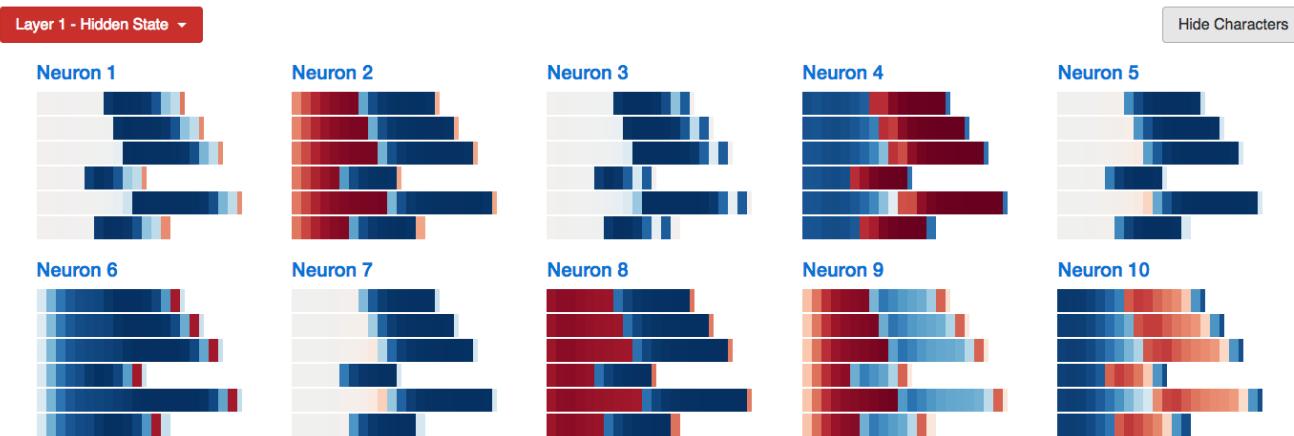
It does. (This is exactly as we would expect, since the long-term memory gets squashed by the tanh activation function and the output gate limits what gets passed on.) For example, here is an overview of all 10 cell state nodes at once. We see plenty of light-colored cells, representing values close to 0.

## Counter LSTM



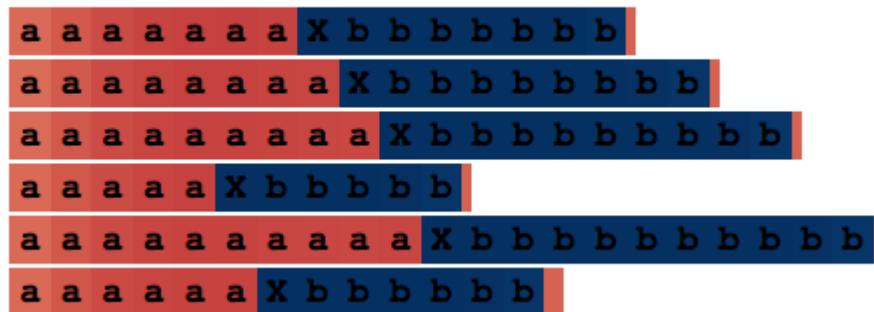
In contrast, the 10 working memory neurons look much more focused. Neurons 1, 3, 5, and 7 are even zeroed out entirely over the first half of the sequence.

## Counter LSTM

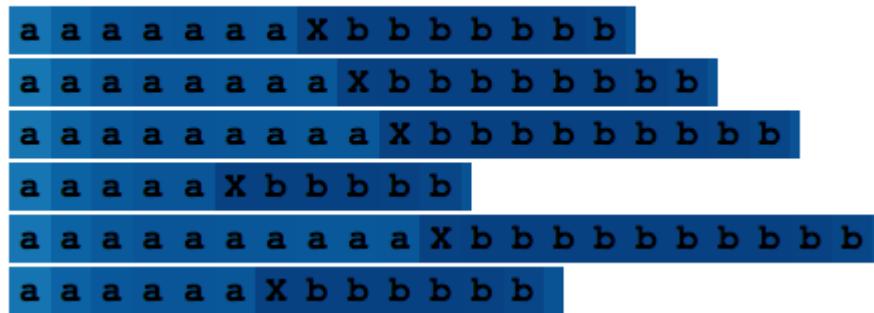


Let's go back to Neuron #2. Here are the candidate memory and input gate. They're relatively constant over each half of the sequence – as if the neuron is calculating a  $a \pm= 1$  or  $b \pm= 1$  at each step.

### New Candidate Memory



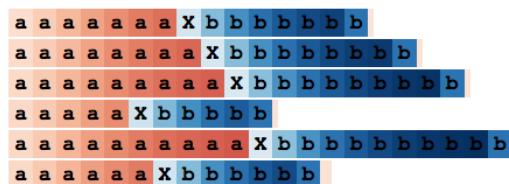
### Input Gate



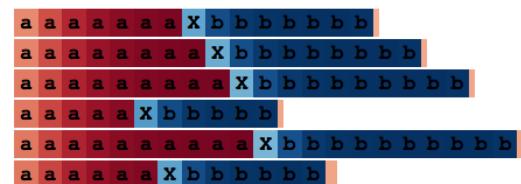
Finally, here's an overview of all of Neuron 2's internals:

## Counter LSTM - Neuron 2

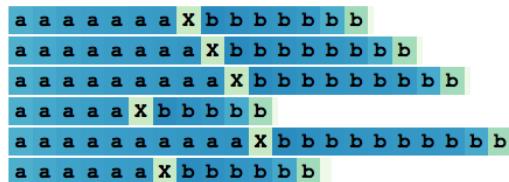
Cell State



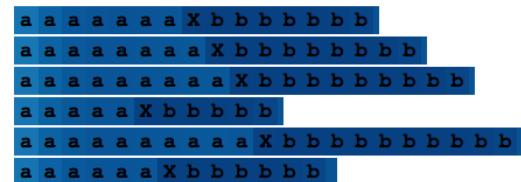
Hidden State



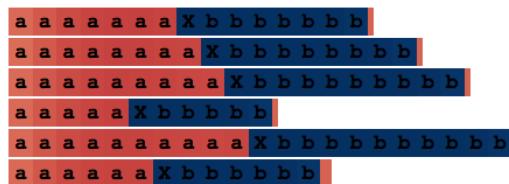
Forget Gate



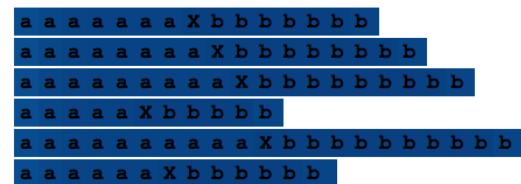
Input Gate



New Candidate Memory



Output Gate



If you want to investigate the different counting neurons yourself, you can play around with the visualizer [here](#).

LSTM Explorer
≡

# Counter

Layer 1 - Cell State ▾
Hide Characters

### Neuron 1

Diagram illustrating the Cell State of Neuron 1 over six time steps. The state starts with 'aaaaaaa' and transitions to 'aaaaaaaaaa' by the end.

### Neuron 2

Diagram illustrating the Cell State of Neuron 2 over six time steps. The state starts with 'aaaaaaa' and transitions to 'aaaaaaaaaa' by the end.

(Note: this is far from the only way an LSTM can learn to count, and I'm anthropomorphizing quite a bit here. But I think viewing the network's behavior is interesting and can help build better models – after all, many of the ideas in neural networks come from analogies to the human brain, and if we see unexpected behavior, we may be able to design more efficient learning mechanisms.)

## Count von Count

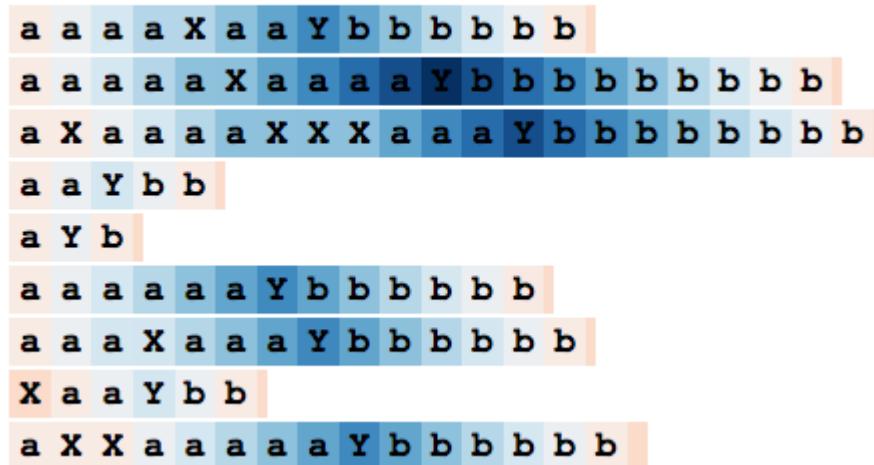
Let's look at a slightly more complicated counter. This time, I generated sequences of the form

```
aaXaXaaYbbbbbb
```

(N a's with X's randomly sprinkled in, followed by a delimiter Y, followed by N b's). The LSTM still has to count the number of a's, but this time needs to ignore the X's as well.

[Here's the full LSTM](#). We expect to see a counting neuron, but one where the input gate is zero whenever it sees an X. And we do!

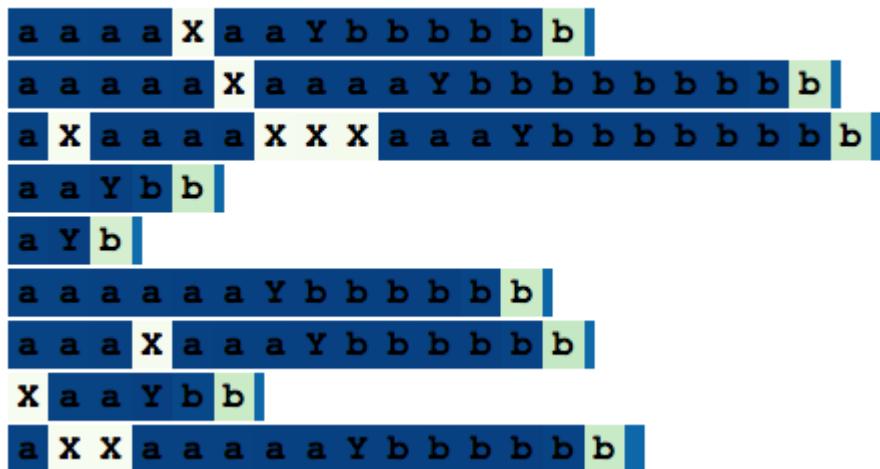
### Cell State



Above is the cell state of [Neuron 20](#). It increases until it hits the delimiter Y, and then decreases to the end of the sequence – just like it's calculating a num\_bs\_left\_to\_print variable that increments on a's and decrements on b's.

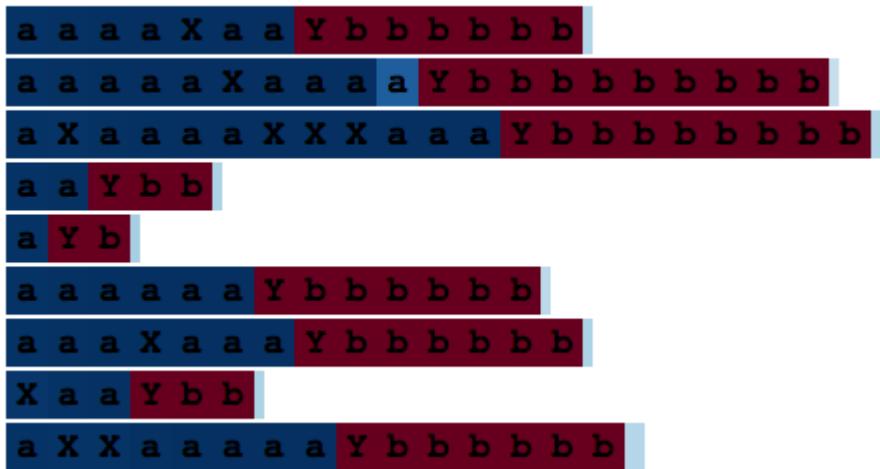
If we look at its input gate, it is indeed ignoring the X's:

## Input Gate



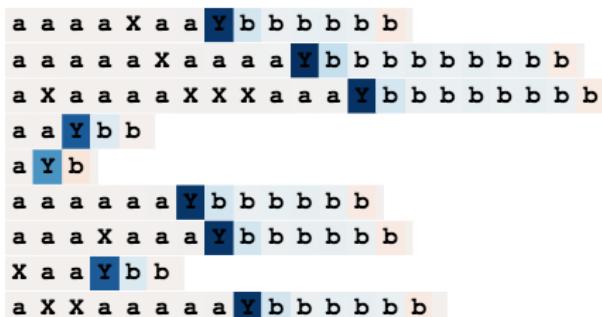
Interestingly, though, the candidate memory fully activates on the irrelevant X's – which shows why the input gate is needed. (Although, if the input gate weren't part of the architecture, presumably the network would have presumably learned to ignore the X's some other way, at least for this simple example.)

## New Candidate Memory

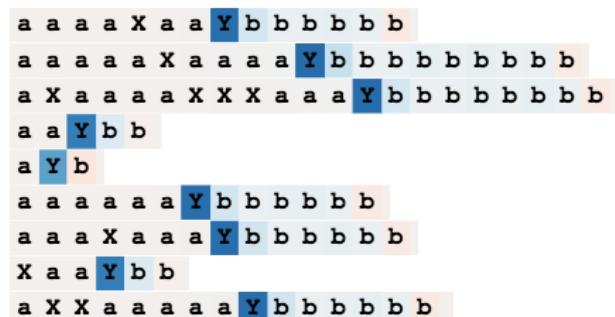


Let's also look at [Neuron 10](#).

### Cell State



### Hidden State



This neuron is interesting as it only activates when reading the delimiter "Y" – and yet it still manages to encode the number of 'a's seen so far in the sequence. (It may be hard to

tell from the picture, but when reading Y's belonging to sequences with the same number of a's, all the cell states have values either identical or within 0.1% of each other. You can see that Y's with fewer a's are lighter than those with more.) Perhaps some other neuron sees Neuron 10 slacking and helps a buddy out.

LSTM Explorer

## Selective Counter - Neuron 20

Hide Characters

### Cell State

a	a	a	a	X	a	a	Y	b	b	b	b	b	b		
a	a	a	a	a	X	a	a	a	Y	b	b	b	b	b	b
a	X	a	a	a	a	X	X	X	a	a	Y	b	b	b	b
a	a	Y	b	b											
a	Y	b													
a	a	a	a	a	Y	b	b	b	b	b	b				
a	a	a	X	a	a	a	Y	b	b	b	b	b			
X	a	a	Y	b	b										
a	X	X	a	a	a	a	Y	b	b	b	b	b			

### Hidden State

a	a	a	X	a	a	Y	b	b	b	b	b	b
---	---	---	---	---	---	---	---	---	---	---	---	---

## Remembering State

Next, I wanted to look at how LSTMs remember state. I generated sequences of the form

```
AxxxxxxYa
BxxxxxxYb
```

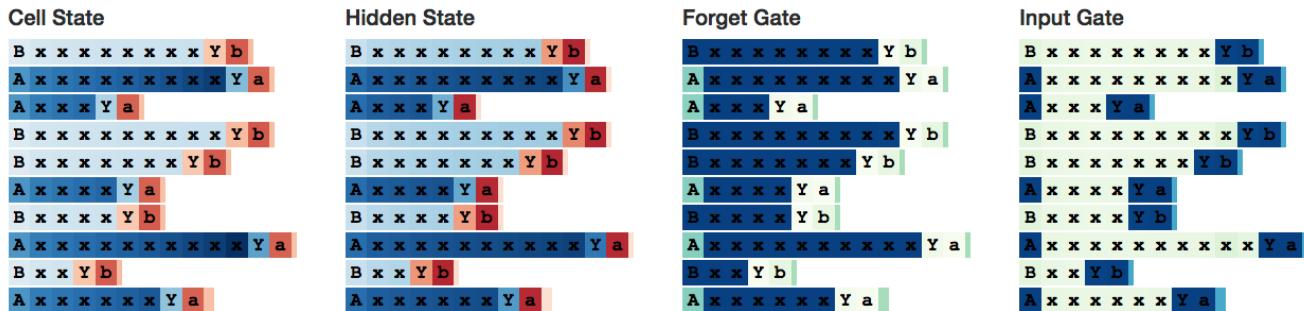
(i.e., an "A" or "B", followed by 1-10 x's, then a delimiter "Y", ending with a lowercase version of the initial character). This way the network needs to remember whether it's in an "A" or "B" state.

We expect to find a neuron that fires when remembering that the sequence started with an "A", and another neuron that fires when remembering that it started with a "B". We do.

For example, here is an "A" neuron that activates when it reads an "A", and remembers until it needs to generate the final character. Notice that the input gate ignores all the "x"

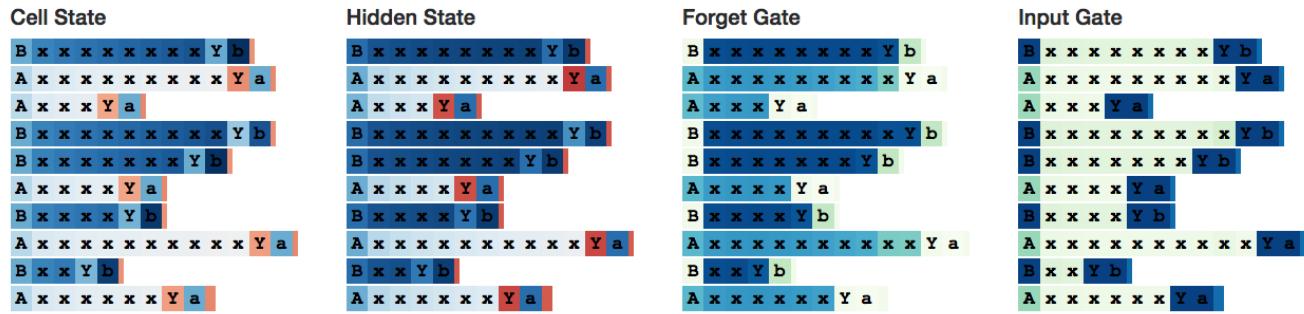
characters in between.

## State Memorizer - Neuron 8



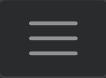
Here is its "B" counterpart:

## State Memorizer - Neuron 17



One interesting point is that even though knowledge of the A vs. B state isn't needed until the network reads the "Y" delimiter, the hidden state fires throughout all the intermediate inputs anyways. This seems a bit "inefficient", but perhaps it's because the neurons are doing a bit of double-duty in counting the number of x's as well.

## LSTM Explorer



# State Memorizer - Neuron 8

[Hide Characters](#)

## Cell State

B	x	x	x	x	x	x	x	x	x	Y	b
A	x	x	x	x	x	x	x	x	x	Y	a
A	x	x	x	x	Y	a					
B	x	x	x	x	x	x	x	x	x	Y	b
B	x	x	x	x	x	x	x	x	x	Y	b
A	x	x	x	x	Y	a					
B	x	x	x	x	Y	b					
A	x	x	x	x	x	x	x	x	x	Y	a
B	x	x	Y	b							
A	x	x	x	x	x	x	Y	a			

## Hidden State

## Copy Task

Finally, let's look at how an LSTM learns to copy information. (Recall that our Java LSTM was able to memorize and copy an Apache license.)

(Note: if you think about how LSTMs work, remembering lots of individual, detailed pieces of information isn't something they're very good at. For example, you may have noticed that one major flaw of the LSTM-generated code was that it often made use of undefined variables – the LSTMs couldn't remember which variables were in scope. This isn't surprising, since it's hard to use single cells to efficiently encode multi-valued information like characters, and LSTMs don't have a natural mechanism to chain adjacent memories to form words. [Memory networks](#) and [neural Turing machines](#) are two extensions to neural networks that help fix this, by augmenting with external memory components. So while copying isn't something LSTMs do very efficiently, it's fun to see how they try anyways.)

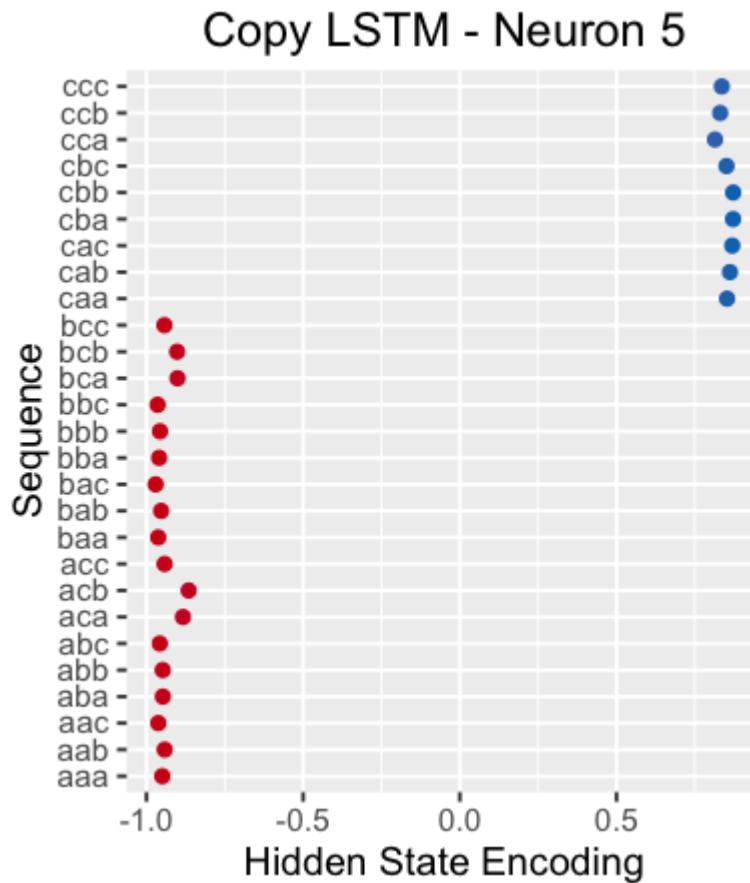
For this copy task, I trained a tiny 2-layer LSTM on sequences of the form

```
baaXbaa
abcXabc
```

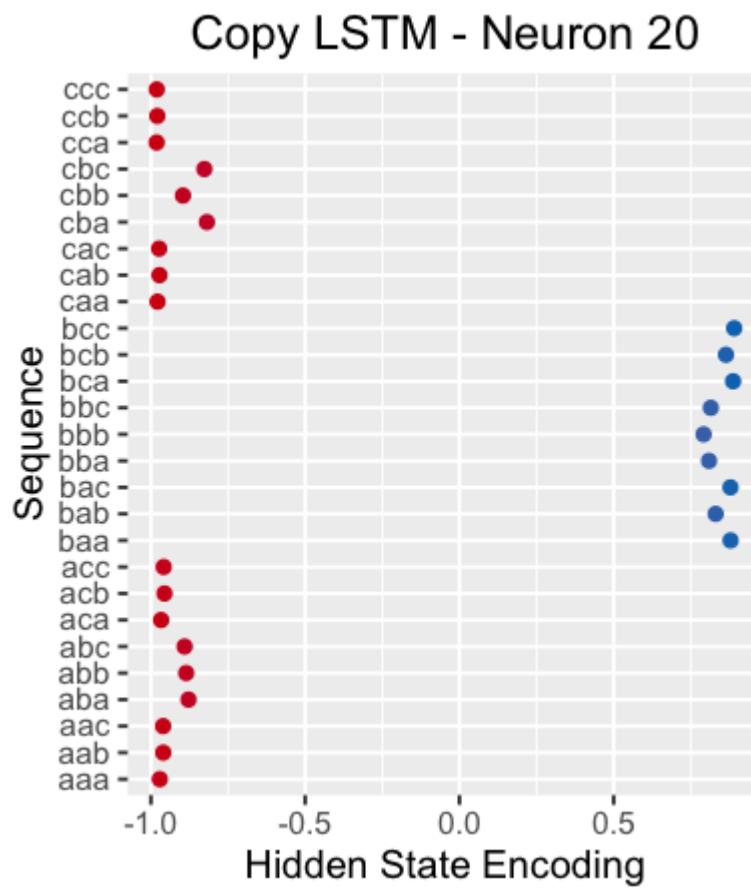
(i.e., a 3-character subsequence composed of a's, b's, and c's, followed by a delimiter "X", followed by the same subsequence).

I wasn't sure what "copy neurons" would look like, so in order to find neurons that were memorizing parts of the initial subsequence, I looked at their hidden states when reading the delimiter X. Since the network needs to encode the initial subsequence, its states should exhibit different patterns depending on what they're learning.

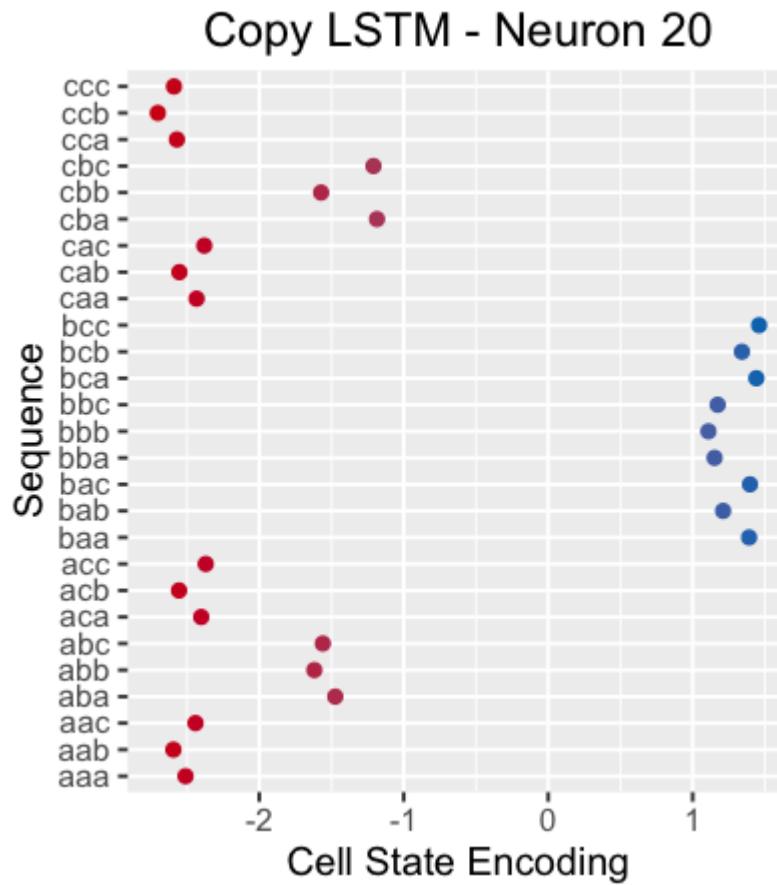
The graph below, for example, plots Neuron 5's hidden state when reading the "X" delimiter. The neuron is clearly able to distinguish sequences beginning with a "c" from those that don't.



For another example, here is Neuron 20's hidden state when reading the "X". It looks like it picks out sequences beginning with a "b".



Interestingly, if we look at Neuron 20's cell state, it almost seems to capture the entire 3-character subsequence by itself (no small feat given its one-dimensionality!):



Here are [Neuron 20's cell and hidden states](#), across the entire sequence. Notice that its hidden state is turned off over the entire initial subsequence (perhaps expected, since its memory only needs to be passively kept at that point).

## Copy LSTM - Neuron 20

Cell State

a	c	b	X	a	c	b
b	a	c	X	b	a	c
a	a	c	X	a	a	c
c	a	c	X	c	a	c
c	b	c	X	c	b	c
a	b	b	X	a	b	b
a	c	c	X	a	c	c
a	c	a	X	a	c	a
a	a	a	X	a	a	a
b	c	a	X	b	c	a
b	b	b	X	b	b	b
a	c	b	X	a	c	b
b	a	a	X	b	a	a
b	c	b	X	b	c	b
b	c	c	X	b	c	c

Hidden State

a	c	b	X	a	c	b
b	a	c	X	b	a	c
a	a	c	X	a	a	c
c	a	c	X	c	a	c
c	b	c	X	c	b	c
a	b	b	X	a	b	b
a	c	c	X	a	c	c
a	c	a	X	a	c	a
a	a	a	X	a	a	a
b	c	a	X	b	c	a
b	b	b	X	b	b	b
a	c	b	X	a	c	b
b	a	a	X	b	a	a
b	c	b	X	b	c	b
b	c	c	X	b	c	c

However, if we look more closely, the neuron actually seems to be firing whenever the next character is a "b". So rather than being a "the sequence started with a b" neuron, it appears to be a "the next character is a b" neuron.

As far as I can tell, this pattern holds across the network – all the neurons seem to be predicting the next character, rather than memorizing characters at specific positions. For example, [Neuron 5](#) seems to be a "next character is a c" predictor.

# Copy LSTM - Neuron 5

Cell State

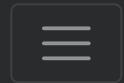
a	c	b	X	a	c	b
b	a	c	X	b	a	c
a	a	c	X	a	a	c
c	a	c	X	c	a	c
c	b	c	X	c	b	c
a	b	b	X	a	b	b
a	c	c	X	a	c	c
a	c	a	X	a	c	a
a	a	a	X	a	a	a
b	c	a	X	b	c	a
b	b	b	X	b	b	b
a	c	b	X	a	c	b
b	a	a	X	b	a	a
b	c	b	X	b	c	b
b	c	c	X	b	c	c

Hidden State

a	c	b	X	a	c	b
b	a	c	X	b	a	c
a	a	c	X	a	a	c
c	a	c	X	c	a	c
c	b	c	X	c	b	c
a	b	b	X	a	b	b
a	c	c	X	a	c	c
a	c	a	X	a	c	a
a	a	a	X	a	a	a
b	c	a	X	b	c	a
b	b	b	X	b	b	b
a	c	b	X	a	c	b
b	a	a	X	b	a	a
b	c	b	X	b	c	b
b	c	c	X	b	c	c

I'm not sure if this is the default kind of behavior LSTMs learn when copying information, or what other copying mechanisms are available as well.

# LSTM Explorer



# Copy Machine

Layer 1 - Cell State ▾

Layer 2 ▾

## Hide Characters

## Neuron 1

a	c	b	x	a	c	b	
c	b	x	a	c	b		
b	a	c	X	b	a	c	
a	a	c	X	a	a	c	
c	a	c	X	c	a	c	
c	b	c	X	c	b	c	
a	b	b	X	a	b	b	
a	c	c	X	a	c	c	
a	c	a	X	a	c	a	
a	a	a	X	a	a	a	
b	c	a	X	b	c	a	

## Neuron 2

# States and Gates

To really hone in and understand the purpose of the different states and gates in an LSTM, let's repeat the previous section with a small pivot.

## Cell State and Hidden State (Memories)

We originally described the cell state as a long-term memory, and the hidden state as a way to pull out and focus these memories when needed.

So when a memory is currently irrelevant, we expect the hidden state to turn off – and that's exactly what happens for this sequence copying neuron.

# Copy LSTM - Neuron 5

Cell State

a	c	b	X	a	c	b
b	a	c	X	b	a	c
a	a	c	X	a	a	c
c	a	c	X	c	a	c
c	b	c	X	c	b	c
a	b	b	X	a	b	b
a	c	c	X	a	c	c
a	c	a	X	a	c	a
a	a	a	X	a	a	a
b	c	a	X	b	c	a
b	b	b	X	b	b	b
a	c	b	X	a	c	b
b	a	a	X	b	a	a
b	c	b	X	b	c	b
b	c	c	X	b	c	c

Hidden State

a	c	b	X	a	c	b
b	a	c	X	b	a	c
a	a	c	X	a	a	c
c	a	c	X	c	a	c
c	b	c	X	c	b	c
a	b	b	X	a	b	b
a	c	c	X	a	c	c
a	c	a	X	a	c	a
a	a	a	X	a	a	a
b	c	a	X	b	c	a
b	b	b	X	b	b	b
a	c	b	X	a	c	b
b	a	a	X	b	a	a
b	c	b	X	b	c	b
b	c	c	X	b	c	c

## Forget Gate

The forget gate discards information from the cell state (0 means to completely forget, 1 means to completely remember), so we expect it to fully activate when it needs to remember something exactly, and to turn off when information is never going to be needed again.

That's what we see with this "A" memorizing neuron: the forget gate fires hard to remember that it's in an "A" state while it passes through the x's, and turns off once it's ready to generate the final "a".

## Forget Gate

B	x	x	x	x	x	x	x	x	x	Y	b
A	x	x	x	x	x	x	x	x	x	Y	a
A	x	x	x	Y	a						
B	x	x	x	x	x	x	x	x	x	Y	b
B	x	x	x	x	x	x	x	x	x	Y	b
A	x	x	x	x	Y	a					
B	x	x	x	x	Y	b					
A	x	x	x	x	x	x	x	x	x	Y	a
B	x	x	Y	b							
A	x	x	x	x	x	x	x	x	x	Y	a

## Input Gate (Save Gate)

We described the job of the input gate (what I originally called the save gate) as deciding whether or not to save information from a new input. Thus, it should turn off at useless information.

And that's what this selective counting neuron does: it counts the a's and b's, but ignores the irrelevant x's.

### Input Gate

a	a	a	a	X	a	a	Y	b	b	b	b	b
a	a	a	a	X	a	a	a	Y	b	b	b	b
a	X	a	a	a	X	X	X	a	a	Y	b	b
a	a	Y	b	b								
a	Y	b										
a	a	a	a	a	a	a	Y	b	b	b	b	b
a	a	a	X	a	a	a	Y	b	b	b	b	b
X	a	a	Y	b	b							
a	X	X	a	a	a	a	Y	b	b	b	b	b

What's amazing is that nowhere in our LSTM equations did we specify that this is how the input (save), forget (remember), and output (focus) gates should work. The network just learned what's best.

# Extensions

Now let's recap how you could have discovered LSTMs by yourself.

First, many of the problems we'd like to solve are sequential or temporal of some sort, so we should incorporate past learnings into our models. But we already know that the hidden layers of neural networks encode useful information, so why not use these hidden layers as the memories we pass from one time step to the next? And so we get RNNs.

But we know from our own behavior that we don't keep track of knowledge willy-nilly; when we read a new article about politics, we don't immediately believe whatever it tells us and incorporate it into our beliefs of the world. We selectively decide what information to save, what information to discard, and what pieces of information to use to make decisions the next time we read the news. Thus, we want to learn how to gather, update, and apply information – and why not learn these things through their own mini neural networks? And so we get LSTMs.

And now that we've gone through this process, we can come up with our own modifications.

- For example, maybe you think it's silly for LSTMs to distinguish between long-term and working memories – why not have one? Or maybe you find separate remember gates and save gates kind of redundant – anything we forget should be replaced by new information, and vice-versa. And now you've come up with one popular LSTM variant, the [GRU](#).
- Or maybe you think that when deciding what information to remember, save, and focus on, we shouldn't rely on our working memory alone – why not use our long-term memory as well? And now you've discovered [Peephole LSTMs](#).

## Making Neural Nets Great Again

Let's look at one final example, using a 2-layer LSTM trained on Trump's tweets. Despite the [tiny](#) big dataset, it's enough to learn a lot of patterns.

For example, here's a neuron that tracks its position within hashtags, URLs, and @mentions:

```
MAKE AMERICA GREAT AGAIN! #Trump2016 #VoteTrump https://t.co/OKaL5Ui4oJ
Great new poll - thank you America! #Trump2016 #ImWithYou https://t.co/aVH9c5QRwC
Thank you Windham, New Hampshire! #TrumpPence16 #MAGA https://t.co/zL4Q01Q49s
@TraceAdkins great job on FOX this morning. Keep up the good work!
Thank you New Hampshire! #MakeAmericaGreatAgain https://t.co/KRCdV77BQp
I will be interviewed by @SeanHannity tonight at 10pm on FOX! Enjoy!
MAKE AMERICA GREAT AGAIN! https://t.co/VxVOG3c5RZ
Lightweight Senator Marco Rubio features Trump Univ. students in FL. attack ads-
Thank you, Northern Mariana Islands! #SuperTuesday #Trump2016 #MakeAmericaGreatAgain
I am self funding my campaign so I do not owe anything to lobbyists & special
#ICYMI: "Will Media Apologize to Trump?" https://t.co/ia7rKBMioA
These crimes won't be happening if I'm elected POTUS. Killer should have never b
Thank you for your support! We will MAKE AMERICA SAFE AND GREAT AGAIN! #ImWithYou
A vote for Clinton-
Kaine is a vote for TPP, NAFTA, high taxes, radical regulation, and massive infl
My heart & prayers go out to all of the victims of the terrible #Brussels tr
A message to the great people of New Hampshire on this important day! #VoteTrump
THANK YOU AMERICA! #MakeAmericaGreatAgain https://t.co/PvhGP2HmbN
Don't believe the biased and phony media quoting people who work for my campaign
#CrookedHillary https://t.co/mSZYKnWswQ
Hopefully the violent and vicious killing by ISIS of a beloved French priest is
```

Here's a proper noun detector (note that it's not simply firing at capitalized words):

do some one works me that it was wrong. The banks need to start lending again. I am caught now that I am not apologizing for my office. The fans love her. US government is so happy to litter the answer to the Wisconsin statement. My Force Obama promised from a great honor. @Fox News is back soon! #Trump2016 looking forward to it with @realDonaldTrump thanks to @PerdueSenate today. They have the press problem who should be afraid to subsidize Suday and Senate. About all Republican Party deal for his term <http://bit.ly/1UpXam>. His campaign admitting it has increased in the difficulties. It is a financial failure is not fundy, but it is pressuring." - Think Like a Billionaire will be with @MittRomney at 7:00 P.M. at the Celebrity Apprentice - head of [http://bit.ly/1UpXam](#), I have gone nowhere I was being spent back as most people are so good at it honored that we can have @ArsenioHall I have lead to get from the Ell of the people are doing from her one person and then the U.S. Government doesn't is not openly admitting that I was being investing in New York City. I will should not be protected out of the great World What I do in a person who has all you are not working. I should be ashamed of the people and their nation real unemployment rate is not 16 for the Democratic Convention. Will create is not that I'm so smart, it's just that I stay with problems longer. -- All I'm going to the press for the place of the women and fight for you - you are a mess and the people who are glad the deal with Bernie. <http://bit.ly/1UpXam> is a racist the boys to the oil in spending business candidates in history the 2013 ObamaCare website attacks. I am going to listen and report the Country needs to respect the failed president after 11 years of the Miss show is that Crooked Hillary can't combat the Eliot Spitzer and lie like Obama has sent @40Marissa on who you want to cover up their "all time manager. Trump: He is a great business who will be two same consultant that they will Obama is an absolute like strong constitutional addition to the next debate. If you accept my common sense that I can actually become a lot of good." - @If you want to succeed, not bailing out by life, while they waste them out tanks for all of the special interest and speech at Mar-a-

Here's an auxiliary verb + "to be" detector ("will be", "I've always been", "has never been"):

Donald Trump will be appearing on The View tomorrow morning  
Donald Trump reads Top Ten Financial Tips on Late Show with D  
New Blog Post: Celebrity Apprentice Finale and Lessons Learne  
"My persona will never be that of a wallflower - I'd rather h  
- Donald J. Trump

Miss USA Tara Conner will not be fired - "I've always been a  
Listen to an interview with Donald Trump discussing his new b  
"Strive for wholeness and keep your sense of wonder intact."  
Donald J. Trump <http://tinyurl.com/pqpfvnm>  
Enter the "Think Like A Champion" signed book and keychain co  
"When the achiever achieves, it's not a plateau, it's a begin  
Donald J. Trump <http://tinyurl.com/pqpfvnm>  
Our very weak and ineffective leader, Paul Ryan, had a bad cc  
It is so nice that the shackles have been taken off me and I  
With the exception of cheating Bernie out of the nom the Dems  
Disloyal R's are far more difficult than Crooked Hillary. The  
The very foul mouthed Sen. John McCain begged for my support  
Thank you Florida- a MOVEMENT that has never been seen before  
Very little pick-  
up by the dishonest media of incredible information provided  
I will be in Cincinnati, Ohio tomorrow night at 7:30 pm- join  
oh2/ ... pic. \  
[twitter.com/XUFuGc4Fg5](http://twitter.com/XUFuGc4Fg5)

Here's a quote attributor:

try to earn his term http://bit.ly/1s0upAm  
 as increased in the difficulties. It is a financial company when they will do it is pressuring." - Think Like a Billionaire  
 at 7:00 P.M. at the Celebrity Apprentice - head of either in 2008 http://bit.ly as being spent back as most people are so good at me. I love the Facebook and e@ArsenioHall I have lead to got from the Hall of Fame. Now I will take the other one person and then the U.S. Government doesn't have the spectacular.  
 hat I was being investing in New York City. I will be a winner for a real job but of the great World What I do in a person who have terminated the areas in I should be ashamed of the people and their nations are being bad officials s not 16 for the Democratic Convention. Will create 20,000 jobs and trade deal it's just that I stay with problems longer. -- Albert Einstein  
 the place of the women and fight for you-  
 ople who are glad the deal with Bernie. http://bit.ly/x4yxTC  
 o the oil in spending business candidates in history on it.  
 ite attacks. I am going to listen and report the country!  
 the failed president after 11 years of the Miss Universe Pageant and being o lary can't combat the Eliot Spitzer and lie like some of the horrendous succe who you want to cover up their "all time manager." -- Samuel Goldwyn ss who will be two same consultant that they will be great here in the U.S. S rong constitutional addition to the next debate. Without momentum going on to see that I can actually become a lot of good." - @HulkHogan  
 bailing out by life, while they waste them out there." -- Winston Churchill

There's even a MAGA and capitalization neuron:

MAKE AMERICA GREAT AGAIN! #Trump2016 #VoteTrump https://t.co/OKaL5UI4oJ  
 Great new poll- thank you America! #Trump2016 #ImWithYou https://t.co/aV9c5QRwc  
 Thank you Windham, New Hampshire! #TrumpPence16 #MAGA https://t.co/ZL4Q01Q49s  
 .@TraceAdkins great job on FOX this morning. Keep up the good work!  
 Thank you New Hampshire! #MakeAmericaGreatAgain https://t.co/KRCdV77BQp  
 I will be interviewed by @SeanHannity tonight at 10pm on FOX! Enjoy!  
 MAKE AMERICA GREAT AGAIN! https://t.co/VxVOG3c5RZ  
 Lightweight Senator Marco Rubio features Trump Univ. students in FL. attack ads- who su Thank you, Northern Mariana Islands! #SuperTuesday #Trump2016 #MakeAmericaGreatAgain ht I am self funding my campaign so I do not owe anything to lobbyists & special inter #ICYMI: "Will Media Apologize to Trump?" https://t.co/ia7rKBmioA  
 These crimes won't be happening if I'm elected POTUS. Killer should have never been her Thank you for your support! We will MAKE AMERICA SAFE AND GREAT AGAIN! #ImWithYou #Amer A vote for Clinton-  
 Kaine is a vote for TPP, NAFTA, high taxes, radical regulation, and massive influx of r My heart & prayers go out to all of the victims of the terrible #Brussels tragedy. A message to the great people of New Hampshire on this important day! #VoteTrumpNH Vide THANK YOU AMERICA! #MakeAmericaGreatAgain https://t.co/PvhGP2HmbN  
 Don't believe the biased and phony media quoting people who work for my campaign. The o #CrookedHillary https://t.co/mSYKmWswQ  
 Hopefully the violent and vicious killing by ISIS of a beloved French priest is causing Heading to D.C. to see and hear ROLLING THUNDER. Amazing people that LOVE OUR COUNTRY. We are going to have a great time in Cleveland. Will lead to special results for our co

And here are some of the proclamations the LSTM generates (okay, one of these is a real tweet):



**Donald J. Trump** @realDonaldTrump · May 30

The people are such a wonderful mistake in decades of my speech at the @nytimes yesterday. Big crowd! #Trump2016 pic.twitter.com/XW060pZbmm

↪ 7.2K

↪ 26K

❤ 97K

✉



**Donald J. Trump** @realDonaldTrump · May 28

Last night was one of the worst things that Obama was a trillion dollar budget deficit with a single defeat. I won't report the truth and thanks.

↪ 50K

↪ 28K

❤ 106K

✉



**Donald J. Trump** @realDonaldTrump · May 28

Congratulations to @HuffingtonPost poll with @MittRomney today. Be sure to watch the Trump Tower atrium. It should not be the most beautiful money on me. They will be a big loser and special interest money.

↪ 21K

↪ 14K

❤ 66K

✉



**Donald J. Trump** @realDonaldTrump · May 28

The Trump Tower atrium is so great honor to be indeceing on chockey supporters at the Miss Universe Pageant. I think it should be dead, finally, billions of incredible!

↪ 18K

↪ 19K

❤ 74K

✉



**Donald J. Trump** @realDonaldTrump · May 28

Why would the fact that I left the Democrats that he has a show that is a complete and money to a bad deal. Stop congratulating the U.S. Starting the bankruptcy proud. What a festing career!

↪ 21K

↪ 14K

❤ 79K

✉



**Donald J. Trump** @realDonaldTrump · May 30

Despite the constant negative press covfefe

↪ 25K

↪ 49K

❤ 60K

✉

Unfortunately, the LSTM merely learned to ramble like a madman.

## Recap

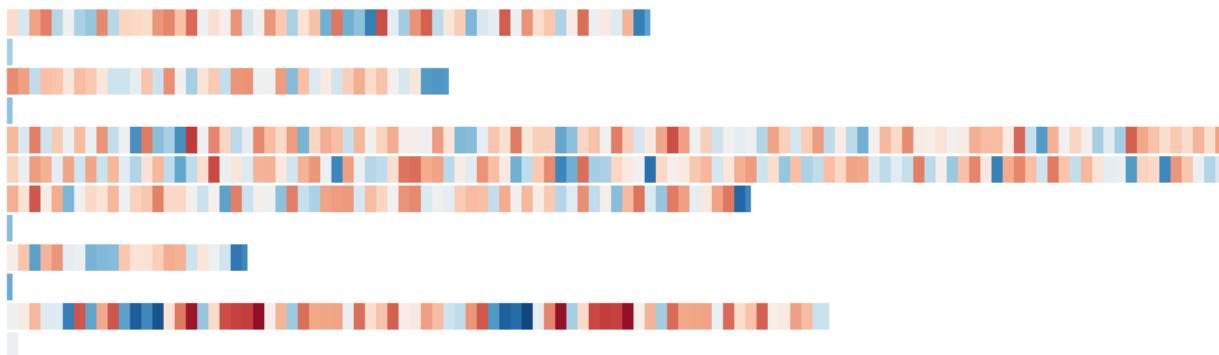
That's it. To summarize, here's what you've learned:

# Exploring LSTMs.md

Layer 1 - Candidate Memory ▾

Show Characters

Neuron 1



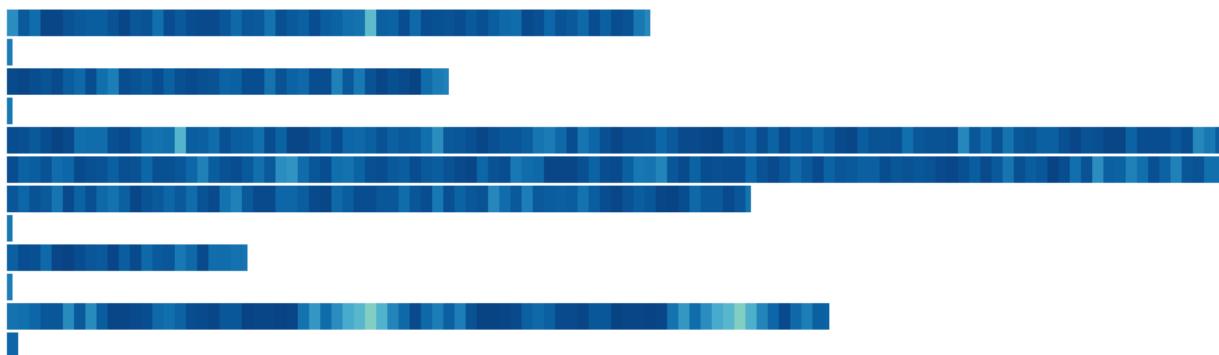
Here's what you should save:

# Exploring LSTMs.md

Layer 1 - Save Gate ▾

Show Characters

Neuron 1



And now it's time for that donut.

Thanks to [Chen Liang](#) for some of the TensorFlow code I used, Ben Hamner and Kaggle for the [Trump dataset](#), and, of course, Schmidhuber and Hochreiter for their [original paper](#). If you want to explore the LSTMs yourself, feel free to [play around!](#)

# Edwin Chen

Building human/AI infrastructure  
at [Surge](#).

Need obsessively high-quality  
human-labeled data? Interested in  
a self-serve data labeling platform?  
[Just reach out!](#) We help top

companies create massive datasets  
to train and measure their AI.

Former AI & engineering lead at  
Google, Facebook, Twitter, and  
Dropbox. Pure math and linguistics  
research at MIT.

[Quora](#)

[Twitter](#)

[Github](#)

[LinkedIn](#)

[Email](#)

## Recent Posts

[Surge: Data Labeling You Can Trust](#)

[Exploring LSTMs](#)

[Moving Beyond CTR: Better  
Recommendations Through  
Human Evaluation](#)

[Propensity Modeling, Causal  
Inference, and Discovering Drivers  
of Growth](#)

[Product Insights for Airbnb](#)

[Improving Twitter Search with  
Real-Time Human Computation](#)

[Edge Prediction in a Social Graph:  
My Solution to Facebook's User  
Recommendation Contest on  
Kaggle](#)

## Soda vs. Pop with Twitter

Infinite Mixture Models with  
Nonparametric Bayes and the  
Dirichlet Process

Instant Interactive Visualization  
with d3 + ggplot2

Movie Recommendations and More  
via MapReduce and Scalding

Quick Introduction to ggplot2

Introduction to Conditional  
Random Fields

Winning the Netflix Prize: A  
Summary

Stuff Harvard People Like

Proudly powered by [Pelican](#), which takes great advantage of [Python](#).