Notebooks =)
① Sentiment Classification - Model - LSTM
② Stock-Price Prediction - LSTM
③ Language Model - RNN | GRU | LSTM

# n-gram Language Models: Example

Suppose we are learning a 4-gram Language Model.

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

$$P(w|\text{students opened their}) = \frac{\text{count(students opened their } w)}{\text{count(students opened their)}}$$

For example, suppose that in the corpus:
- "students opened their" occurred 1000 times
- "students opened their books" occurred 400 times
  - → P(books | students opened their) = 0.4
- "students opened their exams" occurred 100 times
  - → P(exams | students opened their) = 0.1

9

# Sparsity Problems with n-gram Language Models

Sparsity Problem 1

**Problem:** What if "students opened their $w$" never occurred in data? Then $w$ has probability 0!

**(Partial) Solution:** Add small $\delta$ to the count for every $w \in V$. This is called *smoothing*.

$$P(w|\text{students opened their}) = \frac{\text{count(students opened their } w)}{\text{count(students opened their)}}$$

Sparsity Problem 2

**Problem:** What if "students opened their" never occurred in data? Then we can't calculate probability for *any* $w$!

**(Partial) Solution:** Just condition on "opened their" instead. This is called *backoff*.

In general - books
than
exams
but if we know course
than it ex

xt

gm

# A RNN Language Model

$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$

books

laptops

a ... zoo

**output distribution**

$\hat{y}^{(t)} = \text{softmax}\left(Uh^{(t)} + b_2\right) \in \mathbb{R}^{|V|}$

$U$

$h^{(0)}$  $h^{(1)}$  $h^{(2)}$  $h^{(3)}$  $h^{(4)}$

**hidden states**

$h^{(t)} = \sigma\left(W_h h^{(t-1)} + W_e e^{(t)} + b_1\right)$

$h^{(0)}$ is the initial hidden state

$W_h \quad W_h \quad W_h \quad W_h$

$W_e \quad W_e \quad W_e \quad W_e$

**word embeddings**

$e^{(t)} = Ex^{(t)}$

$e^{(1)} \quad e^{(2)} \quad e^{(3)} \quad e^{(4)}$

$E \quad E \quad E \quad E$

**words / one-hot vectors**

$x^{(t)} \in \mathbb{R}^{|V|}$

the      students     opened      their

$x^{(1)}$     $x^{(2)}$     $x^{(3)}$     $x^{(4)}$

23

*Note*: this input sequence could be much longer, but this slide doesn't have space!

---

# A RNN Language Model

$\hat{y}^{(4)} = P(x^{(5)}|\text{the students opened their})$
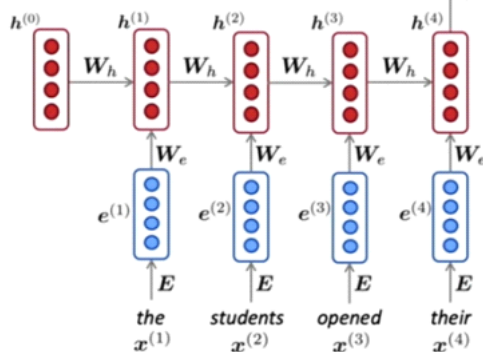
books

laptops

a ... zoo

**RNN Advantages:**
- Can process any length input
- Computation for step *t* can (in theory) use information from many steps back
- Model size doesn't increase for longer input
- Same weights applied on every timestep, so there is symmetry in how inputs are processed.

**RNN Disadvantages:**
- Recurrent computation is slow
- In practice, difficult to access information from many steps back

$h^{(0)}$  $h^{(1)}$  $h^{(2)}$  $h^{(3)}$  $h^{(4)}$

$U$

$W_h \quad W_h \quad W_h \quad W_h$

$W_e \quad W_e \quad W_e \quad W_e$

$e^{(1)} \quad e^{(2)} \quad e^{(3)} \quad e^{(4)}$

$E \quad E \quad E \quad E$

the      students     opened      their

$x^{(1)}$     $x^{(2)}$     $x^{(3)}$     $x^{(4)}$

24

$t = 1, 2$

$t = 1 \_ \_ 1000$

# Training a RNN Language Model

Loss ⟶ $J^{(1)}(\theta)$ + $J^{(2)}(\theta)$ + $J^{(3)}(\theta)$ + $J^{(4)}(\theta)$ + ... = $J(\theta) = \frac{1}{T}\sum_{t=1}^{T} J^{(t)}(\theta)$

Predicted prob dists ⟶ $\hat{y}^{(1)}$  $\hat{y}^{(2)}$  $\hat{y}^{(3)}$  $\hat{y}^{(4)}$

$U$

$h^{(0)}$  $h^{(1)}$  $h^{(2)}$  $h^{(3)}$  $h^{(4)}$

$W_h$

$W_e$

$e^{(1)}$  $e^{(2)}$  $e^{(3)}$  $e^{(4)}$

$E$

Corpus ⟶ the  students  opened  their  exams  ...

$x^{(1)}$  $x^{(2)}$  $x^{(3)}$  $x^{(4)}$

# Generating text with a RNN Language Model

Just like a n-gram Language Model, you can use a RNN Language Model to generate text by repeated sampling. Sampled output is next step's input.

favorite  season  is  spring

sample  sample  sample  sample

$\hat{y}^{(1)}$  $\hat{y}^{(2)}$  $\hat{y}^{(3)}$  $\hat{y}^{(4)}$

$U$

$h^{(0)}$  $h^{(1)}$  $h^{(2)}$  $h^{(3)}$  $h^{(4)}$

$W_h$

$W_e$

$e^{(1)}$  $e^{(2)}$  $e^{(3)}$  $e^{(4)}$

$E$

my  favorite  season  is  spring

# Generating text with a RNN Language Model

- Let's have some fun!
- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on *Harry Potter*:

"Sorry," Harry shouted, panicking—"I'll leave those brooms in London, are they?"

"No idea," said Nearly Headless Nick, casting low close by Cedric, carrying the last bit of treacle Charms, from Harry's shoulder, and to answer him the common room perched upon it, four arms held a shining knob from when the spider hadn't felt it seemed. He reached the teams too.

**Source:** https://medium.com/deep-writing/harry-potter-written-by-artificial-intelligence-8a9431803da6

# Generating text with a RNN Language Model

- Let's have some fun!
- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on recipes:

```
    Title: CHOCOLATE RANCH BARBECUE
Categories: Game, Casseroles, Cookies, Cookies
     Yield: 6 Servings

     2 tb Parmesan cheese -- chopped
     1 c  Coconut milk
     3    Eggs, beaten
```

Place each pasta over layers of lumps. Shape mixture into the moderate oven and simmer until firm. Serve hot in bodied fresh, mustard, orange and cheese.

Combine the cheese and salt together the dough in a large skillet; add the ingredients and stir in the chocolate and pepper.

**Source:** https://gist.github.com/nylki/1efbaa36635956d35bcc

# Generating text with a RNN Language Model

- Let's have some fun!
- You can train a RNN-LM on any kind of text, then generate text in that style.
- RNN-LM trained on paint color names:

| Color | Name | RGB | | Color | Name | RGB |
|---|---|---|---|---|---|---|
| | Ghasty Pink | 231 137 165 | | | Sand Dan | 201 172 143 |
| | Power Gray | 151 124 112 | | | Grade Bat | 48 94 83 |
| | Navel Tan | 199 173 140 | | | Light Of Blast | 175 150 147 |
| | Bock Coe White | 221 215 236 | | | Grass Bat | 176 99 108 |
| | Horble Gray | 178 181 196 | | | Sindis Poop | 204 205 194 |
| | Homestar Brown | 133 104 85 | | | Dope | 219 209 179 |
| | Snader Brown | 144 106 74 | | | Testing | 156 101 106 |
| | Golder Craam | 237 217 177 | | | Stoner Blue | 152 165 159 |
| | Hurky White | 232 223 215 | | | Burble Simp | 226 181 132 |
| | Burf Pink | 223 173 179 | | | Stanky Bean | 197 162 171 |
| | Rose Hork | 230 215 198 | | | Turdly | 190 164 116 |

This is an example of a character-level RNN-LM (predicts what character comes next)

**Source:** http://aiweirdness.com/post/160776374467/new-paint-colors-invented-by-neural-network

# Why should we care about Language Modeling?

- Language Modeling is a benchmark task that helps us measure our progress on understanding language

- Language Modeling is a subcomponent of many NLP tasks, especially those involving generating text or estimating the probability of text:
    - Predictive typing
    - Speech recognition
    - Handwriting recognition
    - Spelling/grammar correction
    - Authorship identification
    - Machine translation
    - Summarization
    - Dialogue
    - etc.

43

# Backpropagation for RNNs



$$\frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h} = \sum_{i=1}^{t} \frac{\partial J^{(t)}}{\partial \boldsymbol{W}_h}\bigg|_{(i)}$$

**Question:** How do we calculate this?

**Answer:** Backpropagate over timesteps *i=t,...,*0, summing gradients as you go.
This algorithm is called
**"backpropagation through time"**

35

# Why is vanishing gradient a problem?
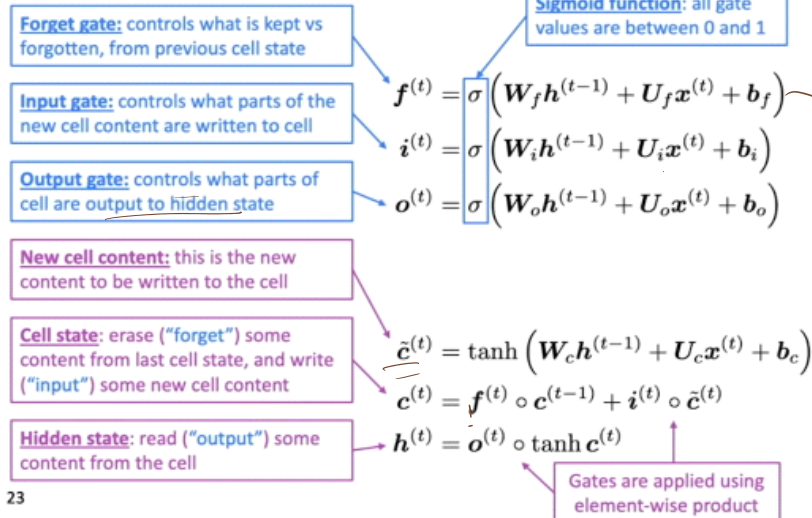
$$J^{(2)}(\theta) \qquad J^{(4)}(\theta)$$

$h^{(1)} \qquad h^{(2)} \qquad h^{(3)} \qquad h^{(4)}$

$W \qquad W \qquad W$

Gradient signal from faraway is lost because it's much smaller than gradient signal from close-by.

So model weights are only updated only with respect to near effects, not long-term effects.
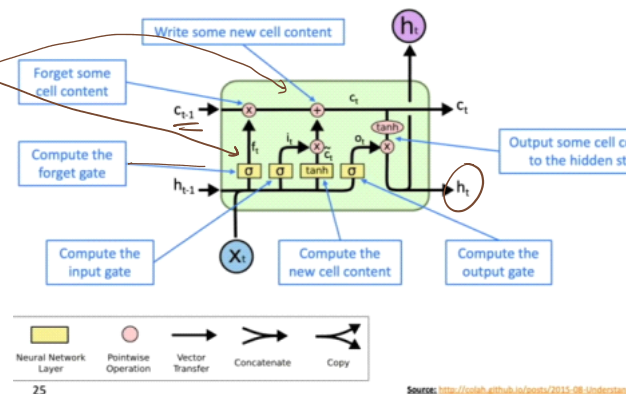
14

# Long Short-Term Memory (LSTM)

We have a sequence of inputs $x^{(t)}$, and we will compute a sequence of hidden states $h^{(t)}$ and cell states $c^{(t)}$. On timestep $t$:

**Forget gate:** controls what is kept vs forgotten, from previous cell state

**Input gate:** controls what parts of the new cell content are written to cell

**Output gate:** controls what parts of cell are output to hidden state

**New cell content:** this is the new content to be written to the cell

**Cell state:** erase ("forget") some content from last cell state, and write ("input") some new cell content

**Hidden state:** read ("output") some content from the cell

23

**Sigmoid function:** all gate values are between 0 and 1

$$f^{(t)} = \sigma\left(W_f h^{(t-1)} + U_f x^{(t)} + b_f\right)$$

$$i^{(t)} = \sigma\left(W_i h^{(t-1)} + U_i x^{(t)} + b_i\right)$$

$$o^{(t)} = \sigma\left(W_o h^{(t-1)} + U_o x^{(t)} + b_o\right)$$

$$\tilde{c}^{(t)} = \tanh\left(W_c h^{(t-1)} + U_c x^{(t)} + b_c\right)$$

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

$$h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

Gates are applied using element-wise product

# Long Short-Term Memory (LSTM)

You can think of the LSTM equations visually like this:



Write some new cell content

Forget some cell content

Compute the forget gate

Output some cell c to the hidden st

Compute the input gate

Compute the new cell content

Compute the output gate

| Neural Network Layer | Pointwise Operation | Vector Transfer | Concatenate | Copy |

25

Source: http://colah.github.io/posts/2015-08-Understan

ontent
ate

ding LSTMs/

# Gated Recurrent Units (GRU)

- Proposed by Cho et al. in 2014 as a simpler alternative to the LSTM.
- On each timestep $t$ we have input $x^{(t)}$ and hidden state $h^{(t)}$ (no cell state).

**Update gate:** controls what parts of hidden state are updated vs preserved

**Reset gate:** controls what parts of previous hidden state are used to compute new content

**New hidden state content:** reset gate selects useful parts of prev hidden state. Use this and current input to compute new hidden content.

**Hidden state:** update gate simultaneously controls what is kept from previous hidden state, and what is updated to new hidden state content

$$u^{(t)} = \sigma\left(W_u h^{(t-1)} + U_u x^{(t)} + b_u\right)$$

$$r^{(t)} = \sigma\left(W_r h^{(t-1)} + U_r x^{(t)} + b_r\right)$$

$$\tilde{h}^{(t)} = \tanh\left(W_h(r^{(t)} \circ h^{(t-1)}) + U_h x^{(t)} + b_h\right)$$

$$h^{(t)} = (1 - u^{(t)}) \circ h^{(t-1)} + u^{(t)} \circ \tilde{h}^{(t)}$$

**How does this solve vanishing gradient?**
Like LSTM, GRU makes it easier to retain info long-term (e.g. by setting update gate to 0)

28    "Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation", Cho et al. 2014, https://arxiv.org/pdf/1406.1078v3.pdf

# LSTM vs GRU

- Researchers have proposed many gated RNN variants, but LSTM and GRU are the most widely-used

- The biggest difference is that GRU is quicker to compute and has fewer parameters

- There is no conclusive evidence that one consistently performs better than the other

- LSTM is a good default choice (especially if your data has particularly long dependencies, or you have lots of training data)

- Rule of thumb: start with LSTM, but switch to GRU if you want something more efficient

29

# Bidirectional RNNs: motivation

Task: Sentiment Classification

**positive**

Sentence encoding

We can regard this hidden state as a representation of the word *"terribly"* in the context of this sentence. We call this a *contextual representation.*

element-wise mean/max        element-wise mean/max

These contextual representations only contain information about the *left* context (e.g. *"the movie was"*).

**What about *right* context?**

In this example, *"exciting"* is in the right context and this modifies the meaning of *"terribly"* (from negative to positive)

the    movie    was    terribly    exciting    !

36

He said, "Teddy Roosevelt was a great President."

He said, "Teddy bears are on sale!"

# Bidirectional RNNs

This contextual representation of *"terribly"* has both left and right context!

Concatenated hidden states

**Backward RNN**

**Forward RNN**

the    movie    was    terribly    exciting    !

37

# Bidirectional RNNs

On timestep $t$:

This is a general notation to mean "compute one forward step of the RNN" – it could be a vanilla, LSTM or GRU computation.

Forward RNN $\quad \overrightarrow{\boldsymbol{h}}^{(t)} = \text{RNN}_{\text{FW}}(\overrightarrow{\boldsymbol{h}}^{(t-1)}, \boldsymbol{x}^{(t)})$

Backward RNN $\quad \overleftarrow{\boldsymbol{h}}^{(t)} = \text{RNN}_{\text{BW}}(\overleftarrow{\boldsymbol{h}}^{(t+1)}, \boldsymbol{x}^{(t)})$

Generally, these two RNNs have separate weights

Concatenated hidden states $\quad \boldsymbol{h}^{(t)} = [\overrightarrow{\boldsymbol{h}}^{(t)}; \overleftarrow{\boldsymbol{h}}^{(t)}]$

We regard this as "the hidden state" of a bidirectional RNN. This is what we pass on to the next parts of the network.
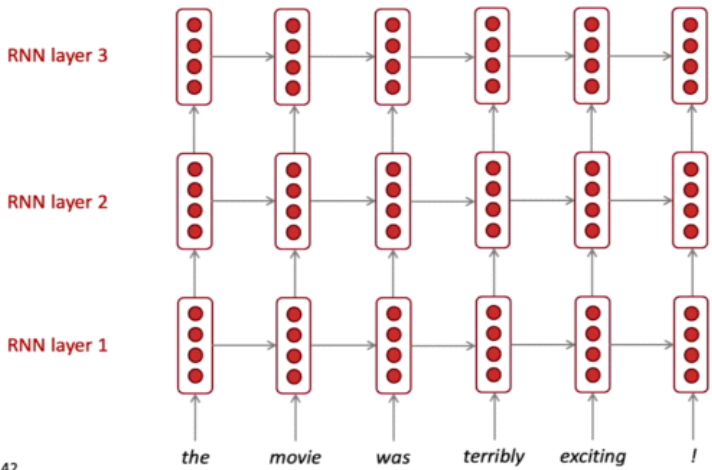
38

# Multi-layer RNNs

The hidden states from RNN layer $i$ are the inputs to RNN layer $i+1$

RNN layer 3

RNN layer 2

RNN layer 1

the    movie    was    terribly    exciting    !

42