

6.036: Introduction to Machine Learning

Final lecture! Thanks for joining us on this adventure!

Lecture start: Tuesdays 9:35am

Who's talking? Prof. Tamara Broderick

Questions? Ask on Piazza: "lecture (week) 13" folder

Materials: slides, video will all be available on Canvas

Live Zoom feed: <https://mit.zoom.us/j/94238622313>

Last Time(s)

- I. State machines & MDPs
- II. Actions change state of world and give reward
- III. Choosing "best" actions

Today's Plan

- I. Back to supervised learning
- II. Sequential data
- III. Recurrent neural networks

Review

Review

expected reward

Review

- $V_{\pi}^h(s)$: expected reward at horizon h with policy π starting at s

Review

- $V_{\pi}^h(s)$: expected reward at horizon h with policy π starting at s
- $Q^h(s, a)$: expected reward of starting at s , making action a , and then making the “best” action for the $h-1$ steps left

Review

- $V_{\pi}^h(s)$: expected reward at horizon h with policy π starting at s
- Can compute value of a policy even if not optimal
- $Q^h(s, a)$: expected reward of starting at s , making action a , and then making the “best” action for the $h-1$ steps left

Review

- $V_{\pi}^h(s)$: expected reward at horizon h with policy π starting at s
 - Can compute value of a policy even if not optimal
- $Q^h(s, a)$: expected reward of starting at s , making action a , and then making the “best” action for the $h-1$ steps left
 - Don't specify a policy π in advance to compute Q . Rather, we typically use Q to figure out the “best” policy.

Review

- $V_{\pi}^h(s)$: expected reward at horizon h with policy π starting at s
 - Can compute value of a policy even if not optimal
- $Q^h(s, a)$: expected reward of starting at s , making action a , and then making the “best” action for the $h-1$ steps left
 - Don't specify a policy π in advance to compute Q .
Rather, we typically use Q to figure out the “best” policy.
- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world

Review

- $V_{\pi}^h(s)$: expected reward at horizon h with policy π starting at s
 - Can compute value of a policy even if not optimal
- $Q^h(s, a)$: expected reward of starting at s , making action a , and then making the “best” action for the $h-1$ steps left
 - Don't specify a policy π in advance to compute Q .
Rather, we typically use Q to figure out the “best” policy.
- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
 - Contrast with *supervised learning*

Review

- $V_{\pi}^h(s)$: expected reward at horizon h with policy π starting at s
 - Can compute value of a policy even if not optimal
- $Q^h(s, a)$: expected reward of starting at s , making action a , and then making the “best” action for the $h-1$ steps left
 - Don’t specify a policy π in advance to compute Q .
Rather, we typically use Q to figure out the “best” policy.
- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
 - Contrast with *supervised learning*
 - Contrast with the Q^h or Q^* function (expected reward of starting at s , making action a , and then making the “best” action ever after)

Review

- $V_{\pi}^h(s)$: expected reward at horizon h with policy π starting at s
 - Can compute value of a policy even if not optimal
- $Q^h(s, a)$: expected reward of starting at s , making action a , and then making the “best” action for the $h-1$ steps left
 - Don’t specify a policy π in advance to compute Q .
Rather, we typically use Q to figure out the “best” policy.
- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
 - Contrast with *supervised learning*
 - Contrast with the Q^h or Q^* function (expected reward of starting at s , making action a , and then making the “best” action ever after)
 - Contrast with (any horizon) *value iteration*

Review

- $V_{\pi}^h(s)$: expected reward at horizon h with policy π starting at s
 - Can compute value of a policy even if not optimal
- $Q^h(s, a)$: expected reward of starting at s , making action a , and then making the “best” action for the $h-1$ steps left
 - Don’t specify a policy π in advance to compute Q .
Rather, we typically use Q to figure out the “best” policy.
- *Reinforcement learning* (RL): learning (to maximize rewards) by interacting with the world
 - Contrast with *supervised learning*
 - Contrast with the Q^h or Q^* function (expected reward of starting at s , making action a , and then making the “best” action ever after)
 - Contrast with (any horizon) *value iteration*
- Today: can use state machines with supervised learning

Text prediction

Text prediction

Final product

VicePresidentOfCompany@HopefullyNotARealEmailA...

Final product

All the documents are finished. Please see attached

Text prediction

Final product

VicePresidentOfCompany@HopefullyNotARealEmailA..

Final product

All the documents are finished. Please see attached

The image shows the Wikipedia homepage with the following language statistics:

Language	Number of Articles
English	6 183 000+
Español	1 637 000+
日本語	1 235 000+
Русский	1 672 000+
Italiano	1 645 000+
Português	1 045 000+
Polski	1 435 000+
Deutsch	2 495 000+
Français	2 262 000+
中文	1 155 000+

The search bar contains the text "autocomp" and shows a dropdown menu with the following results:

- Autocomplete**: Application that predicts the rest of a word a user is typing.
- Search suggest drop-down list**: A list of search suggestions.
- Automotive industry in India**: A travel guide.
- Automotive industry in the United States**: A news source.

Text prediction

Final product

VicePresidentOfCompany@HopefullyNotARealEmailA..

Final product

All the documents are finished. Please see attached



WIKIPEDIA
The Free Encyclopedia

English 6 183 000+ articles	Español 1 637 000+ artículos
日本語 1 235 000+ 記事	Deutsch 2 495 000+ Artikel
Русский 1 672 000+ статей	Français 2 262 000+ articles
Italiano 1 645 000+ voci	中文 1 155 000+ 條目
Português 1 045 000+ artigos	Polski 1 435 000+ haseł

Search bar: autocomp EN [Search Icon]

- Autocomplete**
Application that predicts the rest of a word a user is typing.
- Search suggest drop-down list**
- Automotive industry in India** [ivoyage](#)
e travel guide
- Automotive industry in the United States** [inews](#)
e news source

Text prediction

Final product

VicePresidentOfCompany@HopefullyNotARealEmailA..

Final product

All the documents are finished. Please see attached



WIKIPEDIA

The Free Encyclopedia

English

6 183 000+ articles

Español

1 637 000+ artículos

日本語

1 235 000+ 記事

Deutsch

2 495 000+ Artikel

Русский

1 672 000+ статей

Français

2 262 000+ articles

Italiano

1 645 000+ voci

中文

1 155 000+ 條目



Português

1 045 000+ artigos

Polski

1 435 000+ haseł

autocomp

EN



Autocomplete

Application that predicts the rest of a word a user is typing.

Search suggest drop-down list



Automotive industry in India

[ivoyage](#)
e travel guide



Automotive industry in the United States

[inews](#)
e news source

Text prediction: supervised learning

Text prediction: supervised learning

- Training data: lots of text

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features

label

w

h

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	_
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters.

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we've said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we've said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we’ve said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we’ve said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension
- Idea: just use last character. But lose info

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we’ve said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension
- Idea: just use last character. But lose info
- Idea: use last m characters

Text prediction: supervised learning

- Training data: lots of text
 - “what happens to a dream deferred”

features	label
w	h
wh	a
wha	t
what	–
what_	h
what_h	a
what_ha	p
what_hap	p
what_happ	e

- Classification with 27 classes
- How to featurize?
- Idea: use all previous characters. But so far we’ve said $x^{(i)} \in \mathbb{R}^d$; i.e. fixed dimension
- Idea: just use last character. But lose info
- Idea: use last $m = 3$ characters

A state machine: writing & predicting text

“wha”

A state machine: writing & predicting text

“wha”

A state machine: writing & predicting text

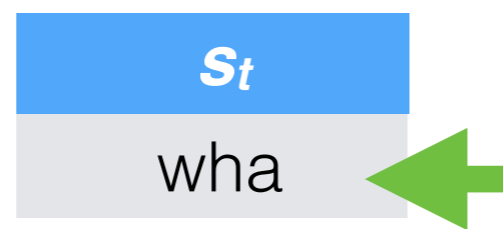
“wha”

A state machine: writing & predicting text



“wha”

A state machine: writing & predicting text



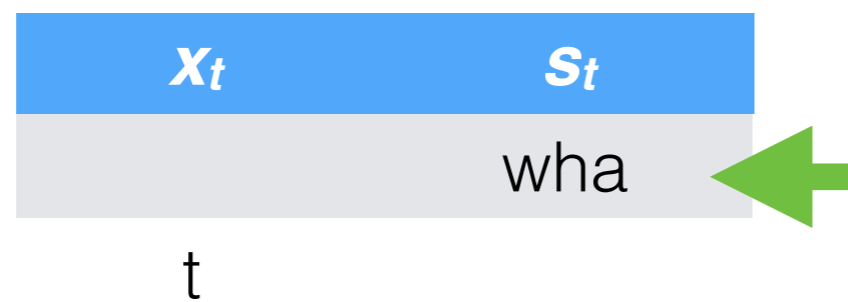
“what”

A state machine: writing & predicting text

“what”

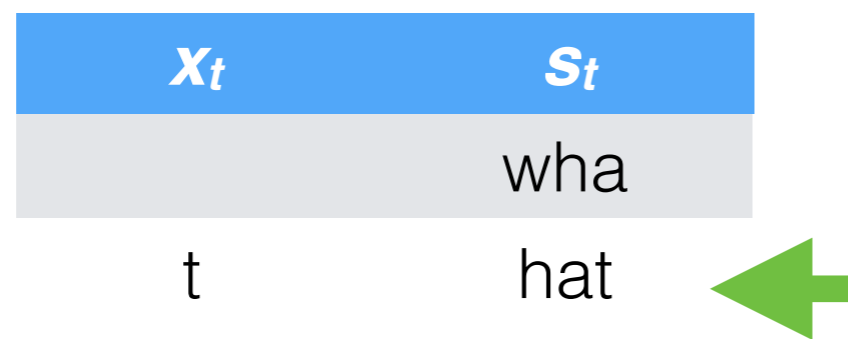


A state machine: writing & predicting text



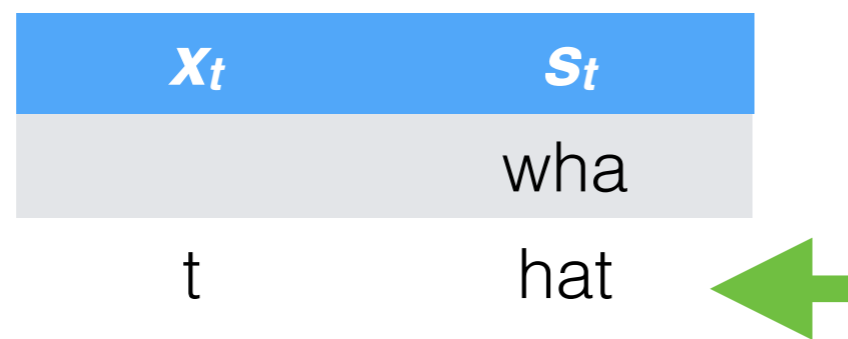
“what”

A state machine: writing & predicting text



“what”

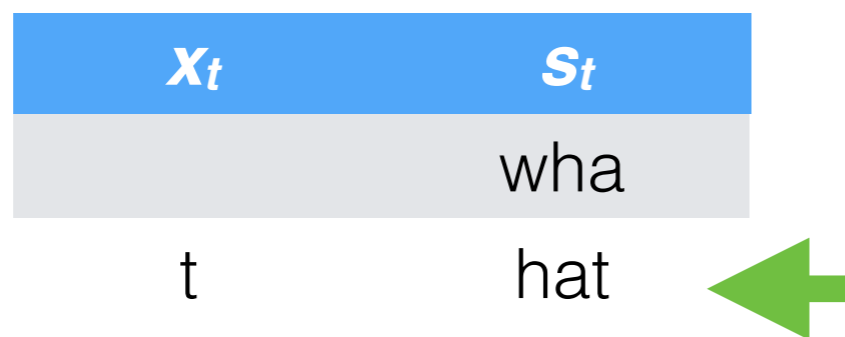
A state machine: writing & predicting text



“what_”

A state machine: writing & predicting text

“what_”



A state machine: writing & predicting text

“what_”

x_t	s_t
	wha
t	hat
_	at_ ←

A state machine: writing & predicting text

“what happens to a
dream deferred”

x_t	s_t
	wha
t	hat
_	at_
h	t_h

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}

“what happens to a dream deferred”

x_t	s_t
	wha
t	hat
_	at_
h	t_h

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
- Example:
 - Every sequence of 3 chars

“what happens to a dream deferred”

x_t	s_t
	wha
t	hat
_	at_
h	t_h

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
- Example:
 - Every sequence of 3 chars

“what happens to a
dream deferred”

X_t	S_t
	wha
t	hat
_	at_
h	t_h

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
- Example:
 - Every sequence of 3 chars
 - Every single character

“what happens to a dream deferred”

X_t	S_t
	wha
t	hat
_	at_
h	t_h

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - Every sequence of 3 chars
 - Every single character

“what happens to a dream deferred”

X_t	S_t
	wha
t	hat
_	at_
h	t_h

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$

“what happens to a dream deferred”

x_t	s_t
	wha
t	hat
_	at_
h	t_h

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$

x_t	s_t
	$\wedge\wedge\wedge$
w	$\wedge\wedge w$
h	$\wedge wh$
a	wha
t	hat
_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$

x_t	s_t
	$\wedge\wedge\wedge$
w	$\wedge\wedge w$
h	$\wedge wh$
a	wha
t	hat
_	at_



“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$

x_t	s_t
	$\wedge\wedge\wedge$
w	$\wedge\wedge w$
h	$\wedge wh$
a	wha
t	hat
_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$

x_t	s_t
	$\wedge\wedge\wedge$
w	$\wedge\wedge w$
h	$\wedge wh$
a	wha
t	hat
_	at_

“what happens to a dream deferred”



A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$

x_t	s_t
	$\wedge\wedge\wedge$
w	$\wedge\wedge w$
h	$\wedge wh$
a	wha
t	hat
_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$

x_t	s_t
	$\wedge\wedge\wedge$
w	$\wedge\wedge w$
h	$\wedge wh$
a	wha
t	hat
_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$

previously “y”

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$
 - Output: 27 prediction probs

previously “y”

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

“what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$

previously “y”

- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$
 - Output: 27 prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, \dots, 0.001, 0.01]$

“what happens to a dream deferred”

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$
 - Output: 27 prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, \dots, 0.001, 0.01]$

previously “y”

t	x_t	s_t
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

- $x^{(1)}$: “what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$
 - Output: 27 prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, \dots, 0.001, 0.01]$

previously “y”

t	$x_t^{(1)}$	$s_t^{(1)}$
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

- $x^{(1)}$: “what happens to a dream deferred”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$
 - Output: 27 prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, \dots, 0.001, 0.01]$

previously “y”

t	$x_t^{(1)}$	$s_t^{(1)}$
0		$\wedge\wedge\wedge$
1	w	$\wedge\wedge w$
2	h	$\wedge wh$
3	a	wha
4	t	hat
5	_	at_

- $x^{(1)}$: “what happens to a dream deferred”
- $x^{(2)}$: “if you can keep your head when all about you”
- $x^{(3)}$: “you may write me down in history”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$
 - Output: 27 prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, \dots, 0.001, 0.01]$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$
 - Output: 27 prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, \dots, 0.001, 0.01]$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$
 - Output: 27 prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, \dots, 0.001, 0.01]$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$
 - Output: 27 prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, \dots, 0.001, 0.01]$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s,x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - $s_0 = \wedge\wedge\wedge$
 - E.g. $f(\wedge\wedge w, h) = \wedge wh$
 - Output: 27 prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, \dots, 0.001, 0.01]$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: 27 prediction probs
 - E.g. $p = g(\hat{\wedge}\hat{\wedge}\hat{\wedge}) = [0.08, 0.02, 0.001, 0.01]$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - S_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: 27 prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, 0.001, 0.01]$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - S_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, 0.001, 0.01]$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - S_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g(\wedge\wedge\wedge) = [0.08, 0.02, 0.001, 0.01]$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - S_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - S_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) =$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) =$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} ? \\ ? \\ ? \end{bmatrix} x_t + \begin{bmatrix} ? \\ ? \\ ? \end{bmatrix} s_{t-1}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} ? \\ ? \\ ? \end{bmatrix} x_t + \begin{bmatrix} ? \\ ? \\ ? \end{bmatrix} s_{t-1}$$

3×1

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{matrix} \text{?} \\ \text{?} \\ \text{?} \end{matrix} x_t + \begin{matrix} \text{?} \\ \text{?} \\ \text{?} \end{matrix} s_{t-1}$$

3x1 1x1

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\begin{matrix} s_t \\ 3 \times 1 \end{matrix} = f(s_{t-1}, x_t) = \begin{matrix} \text{?} \\ \text{?} \\ \text{?} \end{matrix} \begin{matrix} x_t \\ 1 \times 1 \end{matrix} + \begin{matrix} \text{?} \\ \text{?} \\ \text{?} \end{matrix} \begin{matrix} s_{t-1} \\ 3 \times 1 \end{matrix}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{matrix} \text{?} \\ \text{?} \\ \text{?} \end{matrix} x_t + \begin{matrix} \text{?} \\ \text{?} \\ \text{?} \end{matrix} s_{t-1}$$

3x1 1x1 3x1

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{matrix} \text{?} \\ \text{?} \\ \text{?} \end{matrix} x_t + \begin{matrix} \text{?} \\ \text{?} \\ \text{?} \end{matrix} s_{t-1}$$

3x1 1x1 3x1

3x1

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\underset{3 \times 1}{s_t} = f(s_{t-1}, x_t) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}_{3 \times 1} \underset{1 \times 1}{x_t} + \text{?}_{3 \times 1} \underset{3 \times 1}{s_{t-1}}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\underset{3 \times 1}{s_t} = f(s_{t-1}, x_t) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}_{3 \times 1} \underset{1 \times 1}{x_t} + \begin{bmatrix} ? \\ ? \\ ? \end{bmatrix}_{3 \times 1} \underset{3 \times 1}{s_{t-1}}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\underset{3 \times 1}{s_t} = f(s_{t-1}, x_t) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}_{3 \times 1} \underset{1 \times 1}{x_t} + \underset{3 \times 3}{\begin{bmatrix} ? \\ ? \\ ? \end{bmatrix}} \underset{3 \times 1}{s_{t-1}}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\begin{matrix}
 s_t & = & f(s_{t-1}, x_t) & = & \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} & x_t & + & \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} & s_{t-1} \\
 3 \times 1 & & & & 3 \times 1 & 1 \times 1 & & 3 \times 3 & 3 \times 1
 \end{matrix}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\begin{matrix}
 s_t = f(s_{t-1}, x_t) = & \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} & x_t + & \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} & s_{t-1} \\
 \text{3x1} & & \text{1x1} & & \text{3x1} \\
 & \text{3x1} & & \text{3x3} &
 \end{matrix}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\begin{matrix}
 s_t & = & f(s_{t-1}, x_t) & = & \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} & x_t & + & \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} & s_{t-1} \\
 3 \times 1 & & & & 3 \times 1 & 1 \times 1 & & 3 \times 3 & 3 \times 1
 \end{matrix}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\begin{matrix}
 s_t & = & f(s_{t-1}, x_t) & = & \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} & x_t & + & \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} & s_{t-1} \\
 3 \times 1 & & & & 3 \times 1 & 1 \times 1 & & 3 \times 3 & 3 \times 1
 \end{matrix}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\begin{array}{l}
 s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} x_t + \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} s_{t-1} \\
 \begin{array}{ccc}
 3 \times 1 & 1 \times 1 & 3 \times 1 \\
 & & 3 \times 3 \\
 & & & 3 \times 1
 \end{array} \\
 p_t = g(s_t) \\
 =
 \end{array}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}_{3 \times 1} x_t + \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}_{3 \times 3} s_{t-1}$$

1×1 3×1

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\underset{3 \times 1}{s_t} = f(s_{t-1}, x_t) = \underset{3 \times 1}{\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}} \underset{1 \times 1}{x_t} + \underset{3 \times 3}{\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}} \underset{3 \times 1}{s_{t-1}}$$

$$\begin{aligned}
 p_t &= g(s_t) \\
 &= f_2(W^o s_t + W_0^o)
 \end{aligned}$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}_{3 \times 1} x_t + \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}_{3 \times 3} s_{t-1}$$

1×1
 3×1

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\underset{3 \times 1}{s_t} = f(s_{t-1}, x_t) = \underset{3 \times 1}{\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}} \underset{1 \times 1}{x_t} + \underset{3 \times 3}{\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}} \underset{3 \times 1}{s_{t-1}}$$

$$\begin{aligned}
 p_t &= g(s_t) \\
 &= f_2(W^o s_t + W_0^o)
 \end{aligned}$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$\underset{3 \times 1}{s_t} = f(s_{t-1}, x_t) = \underset{3 \times 1}{\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}} \underset{1 \times 1}{x_t} + \underset{3 \times 3}{\begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}} \underset{3 \times 1}{s_{t-1}}$$

$$\begin{aligned}
 p_t &= g(s_t) \\
 &= f_2(W^o s_t + W_0^o)
 \end{aligned}$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} x_t + \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} s_{t-1}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} x_t + \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} s_{t-1}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} x_t + \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} s_{t-1}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f(s_{t-1}, x_t) = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1}$$

$$p_t = g(s_t) \\ = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = \begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix}$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f_1 \left(\begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix} \right)$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function f
 - Set of possible outputs \mathcal{Y}
 - Output function $g(s) = g(s)$
- Example: alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

componentwise;
activation
functions

$$s_t = f_1 \left(\begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix} \right)$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f_1 \left(\begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix} \right)$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f_1 \left(\begin{bmatrix} w_1^{sx} \\ w_2^{sx} \\ w_3^{sx} \end{bmatrix} x_t + \begin{bmatrix} w_{11}^{ss} & w_{12}^{ss} & w_{13}^{ss} \\ w_{21}^{ss} & w_{22}^{ss} & w_{23}^{ss} \\ w_{31}^{ss} & w_{32}^{ss} & w_{33}^{ss} \end{bmatrix} s_{t-1} + \begin{bmatrix} w_{0,1}^{ss} \\ w_{0,2}^{ss} \\ w_{0,3}^{ss} \end{bmatrix} \right)$$

$$p_t = g(s_t) \\ = f_2(W^o s_t + W_0^o)$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

same as final layer in all of our neural networks so far, except for notation

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$\begin{aligned} p_t &= g(s_t) \\ &= f_2(W^o s_t + W_0^o) \end{aligned}$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = g(s_t) \\ = f_2(W^o s_t + W_0^o)$$

Put it all together:

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

Put it all together:

$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible inputs \mathcal{X}
 - Initial state
 - Transition function $f(s, x)$
 - Set of possible outputs
 - Output function $p = g(s)$
- Example:
 - Every sequence of 3 chars
 - Every single character
 - s_0
 - E.g. $f\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}, x\right) = \begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$
 - Output: v prediction probs
 - E.g. $p = g\left(\begin{bmatrix} c_1 \\ c_2 \\ c_3 \end{bmatrix}\right)$
- Example with alphabet $\{0, 1\}$ ($v=2$); state is last $m = 3$ chars

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = g(s_t) = f_2(W^o s_t + W_0^o)$$

Put it all together:

$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

“recurrent neural network”

A state machine: writing & predicting text

- Recall state machines:
 - Set of possible states \mathcal{S}
 - Set of possible transitions
 - Initial state
 - Transition function
 - Set of possible outputs
 - Output function
- Example:
 - Every sequence of 3 chars
 - Character

Recall:

- In 2-layer, fully-connected NNs, we learned the features.
- In CNNs, we learned the filters.
- Analogous idea here.

$\begin{bmatrix} c_2 \\ c_3 \\ x \end{bmatrix}$

 transition probs

$$s_t = f_1 (W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss})$$

$$\begin{aligned}
 p_t &= g(s_t) \\
 &= f_2(W^o s_t + W_0^o)
 \end{aligned}$$

Put it all together:

$$p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$$

“recurrent neural network”

Recurrent Neural Networks

Recurrent Neural Networks

Recall: familiar pattern

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

In our example application:

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

In our example application:

1. Use an RNN: $p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

In our example application:

1. Use an RNN: $p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$
$$s_t = f_1 (W^{sx} x_t + W^{ss} s_{t-1} + W_0^{ss})$$

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

In our example application:

1. Use an RNN: $p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$
$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$
$$p_t = f_2(W^o s_t + W_0^o)$$

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

In our example application:

1. Use an RNN: $p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = f_2(W^o s_t + W_0^o)$$

2. Loss: for q sequences, where the i th has length $n^{(i)}$

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

In our example application:

1. Use an RNN: $p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = f_2(W^o s_t + W_0^o)$$

2. Loss: for q sequences, where the i th has length $n^{(i)}$

$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

In our example application:

1. Use an RNN: $p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = f_2(W^o s_t + W_0^o)$$

2. Loss: for q sequences, where the i th has length $n^{(i)}$

$$L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)})$$

element

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

In our example application:

1. Use an RNN: $p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = f_2(W^o s_t + W_0^o)$$

2. Loss: for q sequences, where the i th has length $n^{(i)}$

$$\text{sequence } L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \text{ element}$$

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

In our example application:

1. Use an RNN: $p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = f_2(W^o s_t + W_0^o)$$

2. Loss: for q sequences, where the i th has length $n^{(i)}$

$$\text{sequence } L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \text{ element}$$

$$J(W, W_0) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

Recurrent Neural Networks

Recall: familiar pattern

1. Choose how to predict label (given features & parameters)
2. Choose a loss (between guess & actual label)
3. Choose parameters by trying to minimize the training loss

In our example application:

1. Use an RNN: $p^{(i)} = \text{RNN}(x^{(i)}; W, W_0)$

$$s_t = f_1(W^{sx}x_t + W^{ss}s_{t-1} + W_0^{ss})$$

$$p_t = f_2(W^o s_t + W_0^o)$$

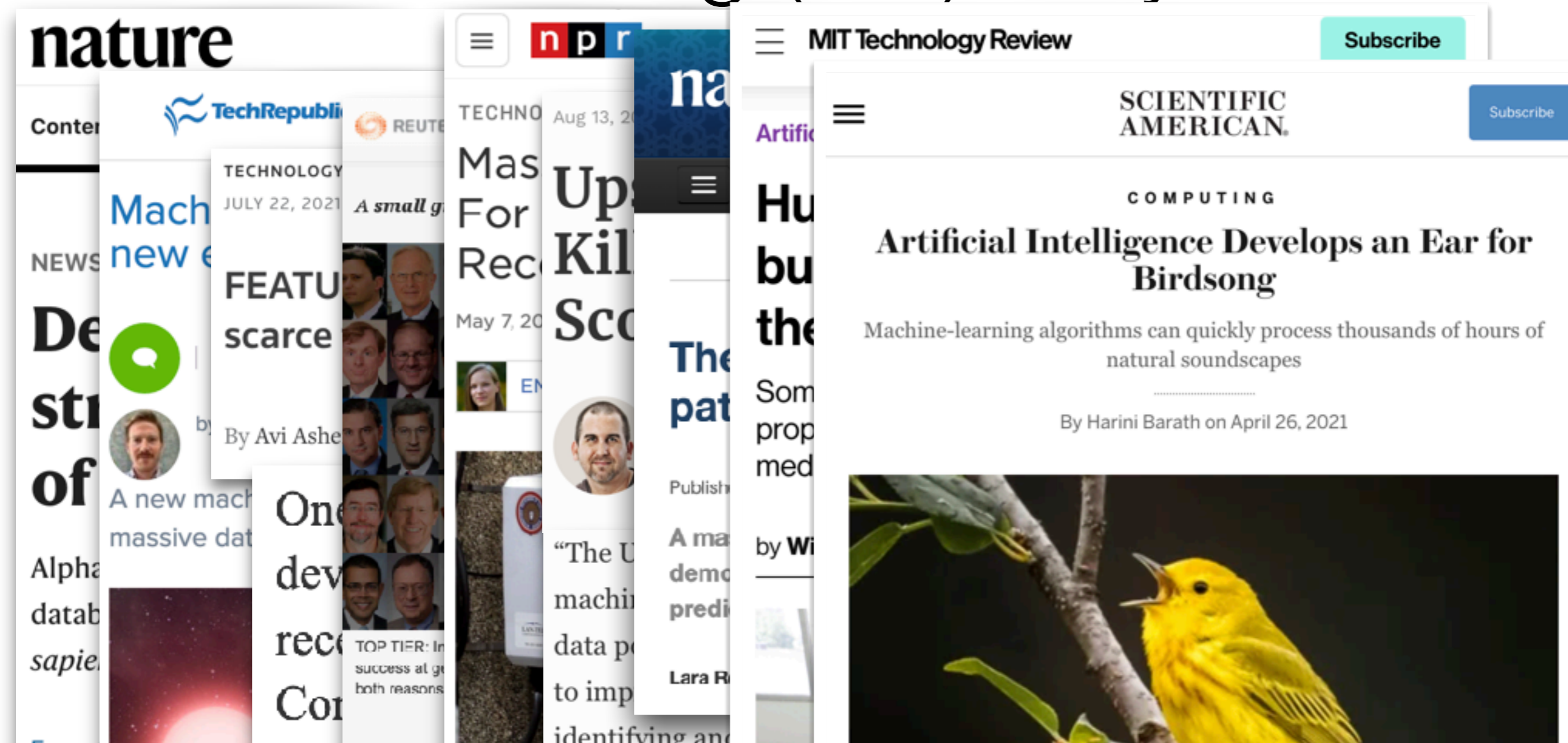
2. Loss: for q sequences, where the i th has length $n^{(i)}$

$$\text{sequence } L_{\text{seq}}(p^{(i)}, y^{(i)}) = \sum_{t=1}^{n^{(i)}} L_{\text{elt}}(p_t^{(i)}, y_t^{(i)}) \text{ element}$$

$$J(W, W_0) = \sum_{i=1}^q L_{\text{seq}}(p^{(i)}, y^{(i)})$$

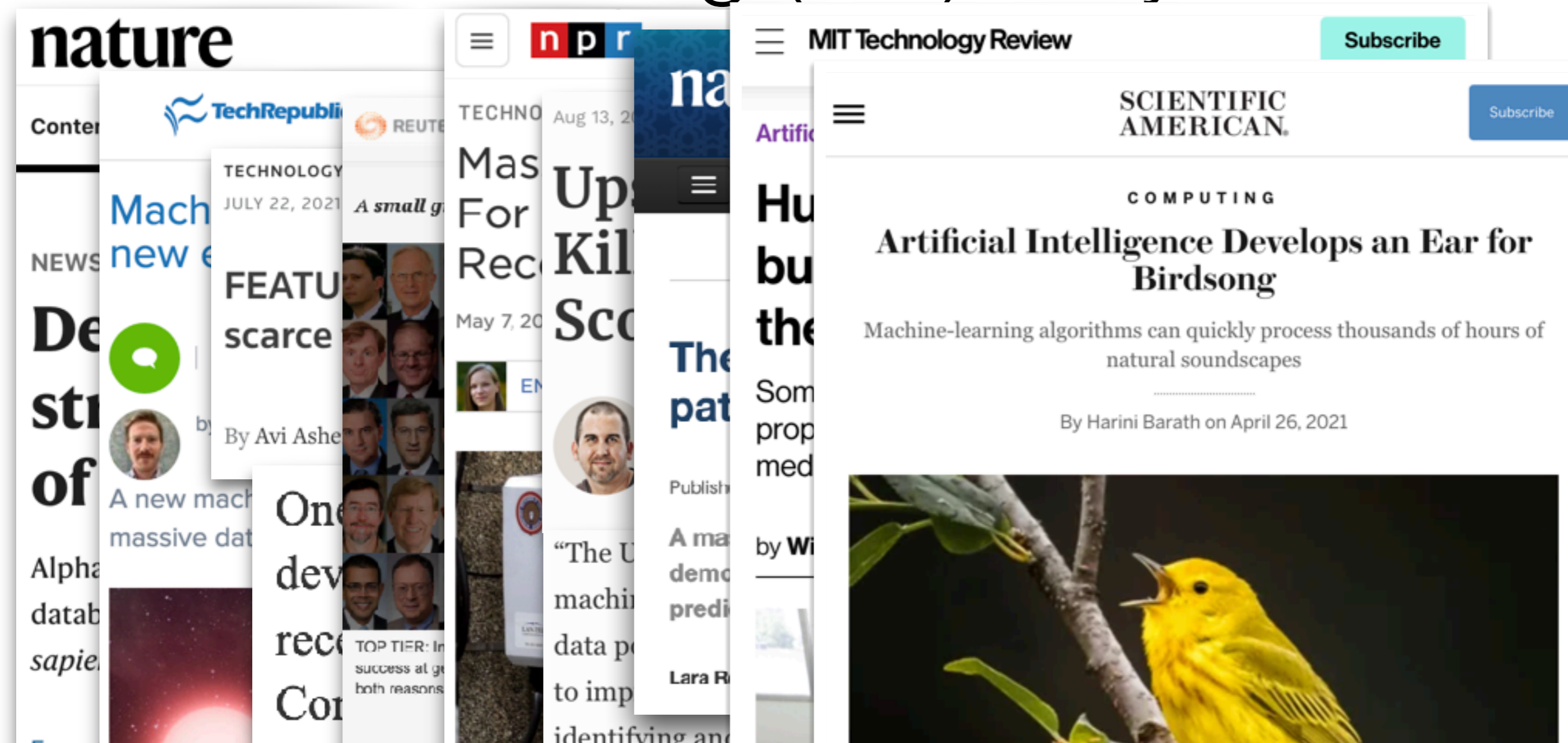
3. Stochastic gradient descent

Machine learning (ML): why & what



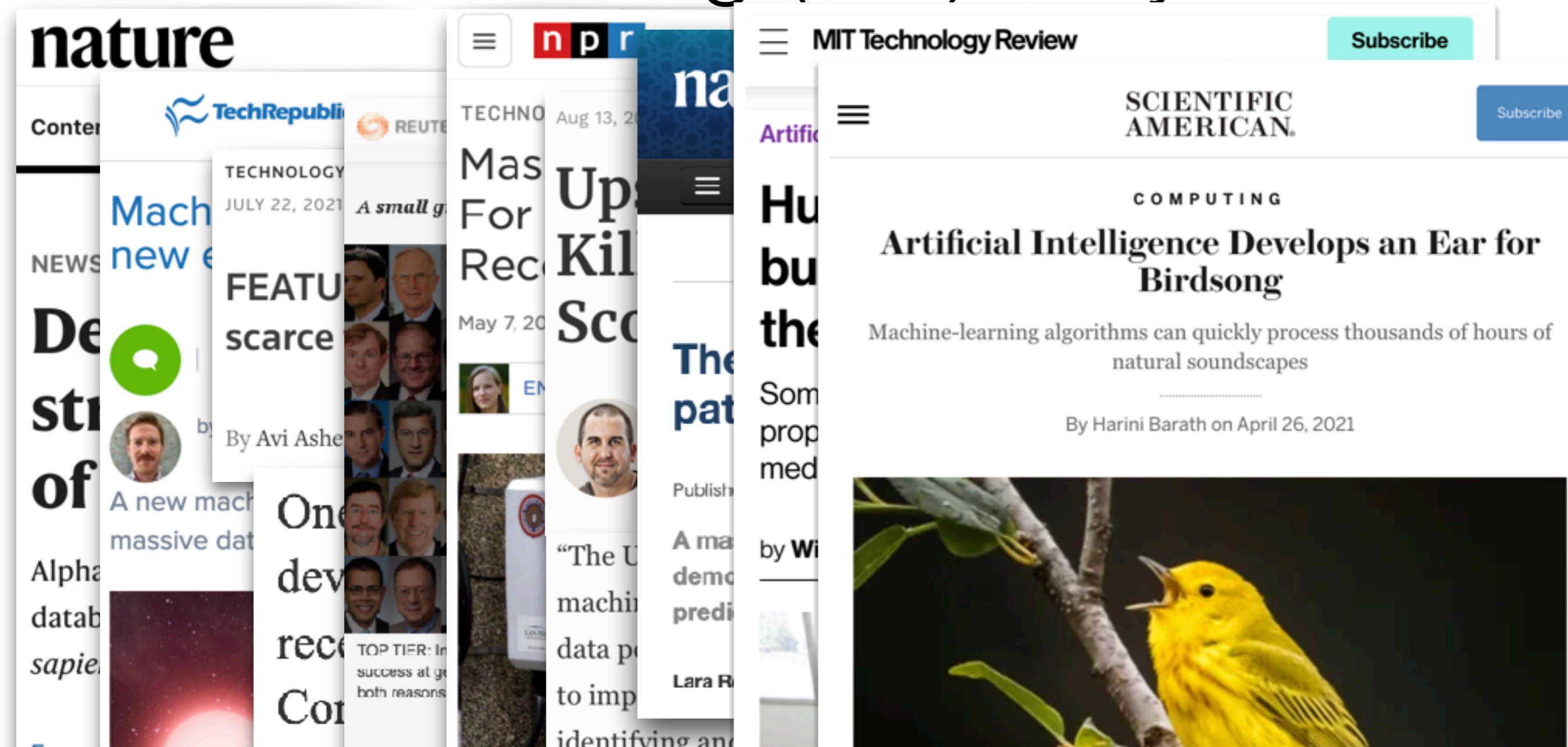
- **What is ML?**

Machine learning (ML): why & what



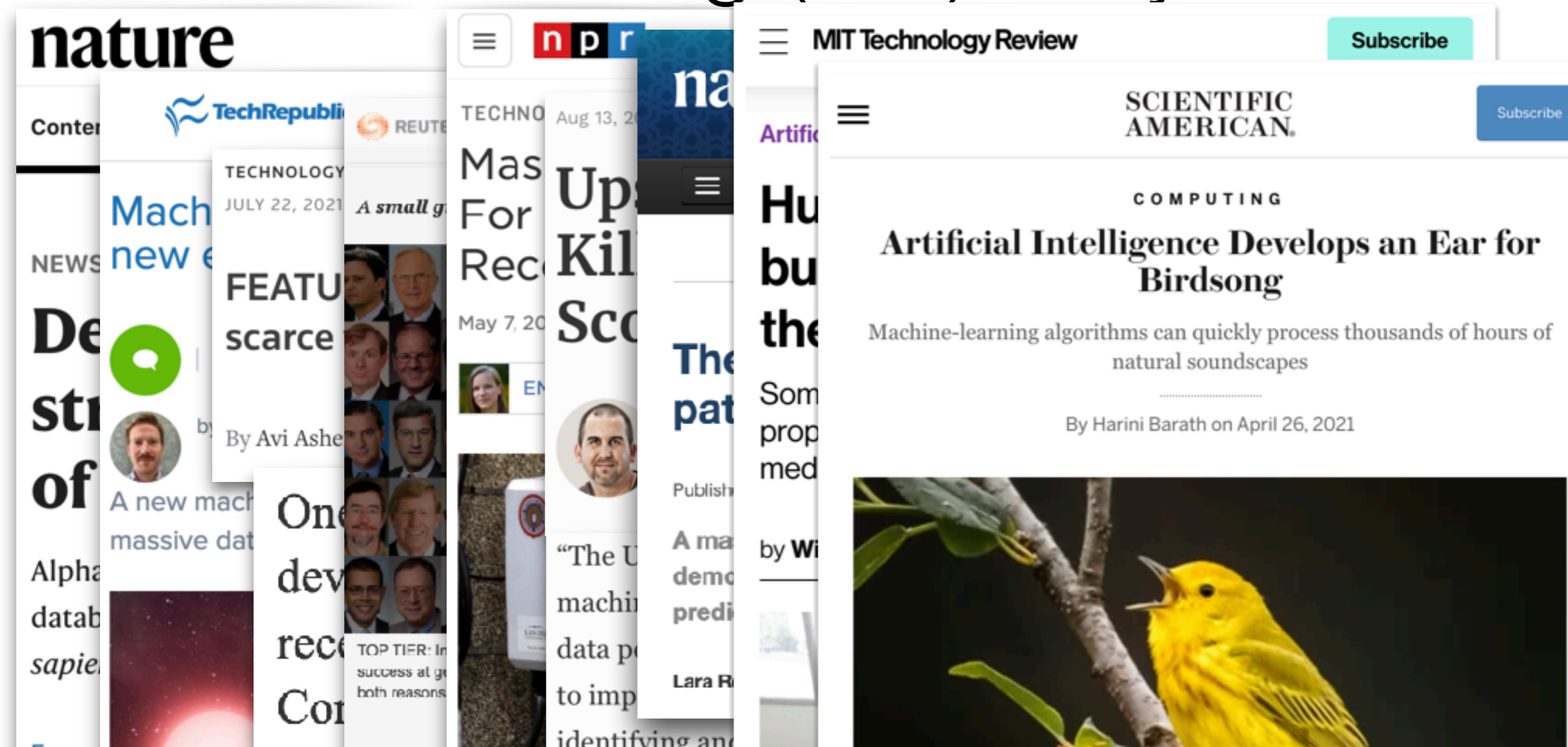
- **What is ML?** A set of methods for making decisions from data. (See the rest of the course!)

Machine learning (ML): why & what



- **What is ML?** A set of methods for making decisions from data. (See the rest of the course!)
- **Why study ML?** To apply; to understand; to evaluate

Machine learning (ML): why & what



- **What is ML?** A set of methods for making decisions from data. (See the rest of the course!)
- **Why study ML?** To apply; to understand; to evaluate
- **Notes:** ML is a tool with pros & cons. ML is built on math



ELECTRICAL ENGINEERING
AND COMPUTER SCIENCE

6.036: Staff

Big thanks to all of our
staff and to the MIT AV
team!

Instructors:



Jehangir
Amjad

Tamara
Broderick

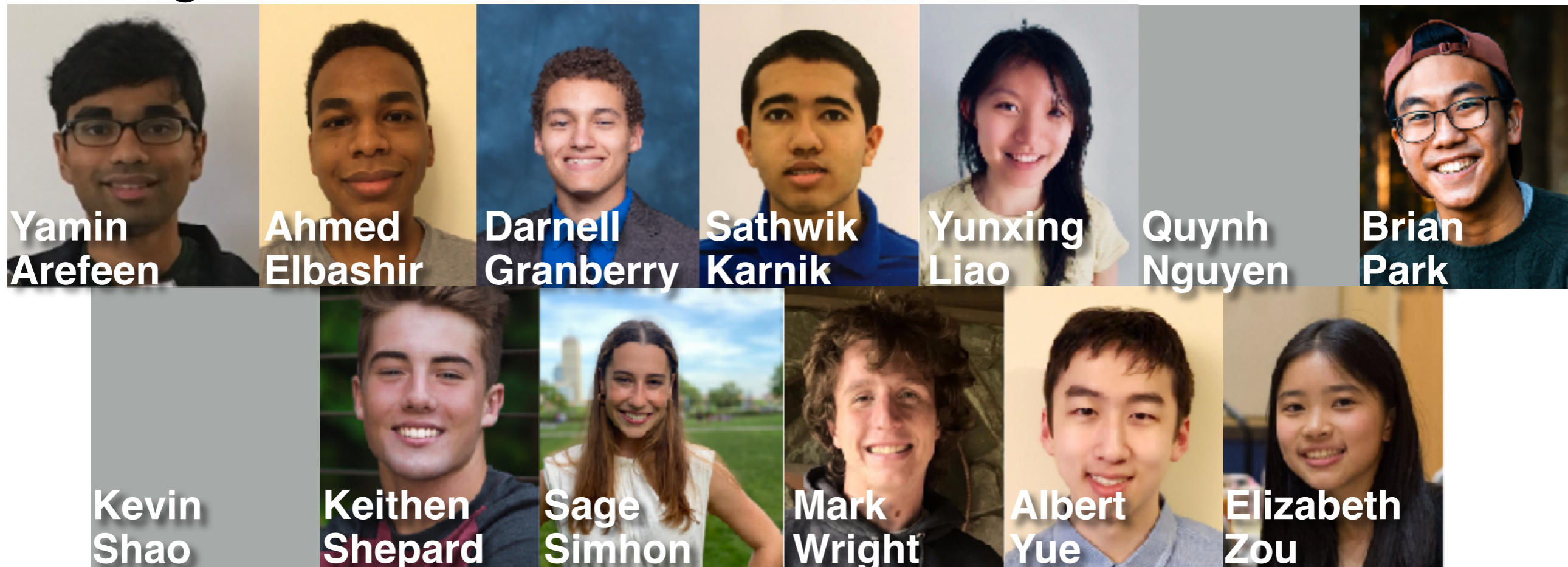
Ike
Chuang

Iddo
Drori

David
Sontag

Tess
Smidt

Teaching Assistants:



Yamin
Arefeen

Ahmed
Elbashir

Darnell
Granberry

Sathwik
Karnik

Yunxing
Liao

Quynh
Nguyen

Brian
Park

Kevin
Shao

Keithen
Shepard

Sage
Simhon

Mark
Wright

Albert
Yue

Elizabeth
Zou

And Lab Assistants!

Thank you!