## Fairness in Reinforcement Learning

## Fairness in Reinforcement Learning

#### Joint work with



Shahin Jabbari



Michael Kearns





Jamie Morgenstern

Aaron Roth



#### Machine learning has 'made it'

#### Machine learning has 'made it'



#### Machine learning has 'made it'



#### Machine learning has 'made it'



#### Machine learning has 'made it'



#### Machine learning has 'made it'

MIT Technology Review

#### An Al-Fueled Credit Formula Might Help You Get a Loan

Startup ZestFinance says it has built a machine-learning system that's smart enough to find new borrowers and keep bias out of its credit analysis.

#### Machine learning has 'made it'

MIT Technology Review

#### An Al-Fueled Credit Formula wight Help You Get a Loan

Startup ZestFinance says it has built a machine-learning system that's smart enough to find new borrowers and keep bias out of its credit analysis.

#### Machine learning has 'made it'

MIT Technology Review

#### An Al-Fueled Credit Formula wight Help You Get a Loan

Startup ZestFinance says it has built a machin system that's smart enough to find new borrov bias out of its credit analysis.



#### Responsibility lagging behind power

#### Responsibility lagging behind power

#### MIT Technology Review

Intelligent Machines

How to Fix Silicon Valley's Sexist Algorithms

Computers are inheriting gender bias implanted in language data sets—and not everyone thinks we should correct it.

#### Responsibility lagging behind power



Computers are inheriting gender bias implanted in language data sets — and not everyone thinks we should correct it.

#### Responsibility lagging behind power

#### PRO PUBLICA Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks. 

Revie Google Image search for CEO Intelligen has Barbie as first female result How to Fix Silicon Valley's Sexist Algorithms

Computers are inheriting gender bias implanted in language data sets - and not everyone thinks we should correct it.

MIT

#### Responsibility lagging behind power

PRO PUBLICA Machine Bias There's soft MIT DEPENDEN'I Techr В **GLE'S ALGORITHM SHOWS Revie** Google How to Fix Si PRESTIGIOUS JOB ADS TO MEN, **Algorithms BUT NOT TO WOMEN** Computers are inheriting g data sets — and not everyone thinks we should correct it

#### Machine learning deserves wide use

#### Machine learning deserves wide use

#### More use means higher stakes

#### Machine learning deserves wide use

#### More use means higher stakes

#### How to balance?

#### **Our Motivation**

#### **Our Motivation**

## Study how machine learning can be efficient, accurate, and fair

## But what is 'fairness'?

## But what is 'fairness'?

# Many reasonable settings and definitions!



1. Specify a setting

1. Specify a setting

#### 2. Specify a fairness definition

1. Specify a setting

#### 2. Specify a fairness definition

#### 3. See what guarantees we can make

1. Specify a setting

#### 2. Specify a fairness definition

3. See what guarantees we can make

## **A General Setting**

## **A General Setting**

#### Learning through feedback from sequential choices

## An Example
### Learning how to grant loans

### Learning how to grant loans



### Learning how to grant loans





### Learning how to grant loans





# \$?

#### Learning how to grant loans



#### Learning how to grant loans



#### Learning how to grant loans



#### Learning how to grant loans



#### Learning how to grant loans





### Fairness (TBD!) adds cost

### Fairness (TBD!) adds cost (performance lower bound)...

### Fairness (TBD!) adds cost (performance lower bound)...

...but cost often not too bad

Fairness (TBD!) adds cost (performance lower bound)...

...but cost often not too bad (performance upper bound)

# Bandits setting is limited – choices don't affect environment

# Bandits setting is limited – choices don't affect environment

\$7



# Bandits setting is limited – choices don't affect environment



\$?

# Bandits setting is limited – choices don't affect environment



# Bandits setting is limited – choices don't affect environment



# Bandits setting is limited – choices don't affect environment

# **Our Specific Setting**

# **Our Specific Setting**



Want to reason about how learning choices affect environment

# **Our Specific Setting**



Want to reason about how learning choices affect environment → reinforcement learning (MDPs)

# **Outline For This Talk**

- Specify a setting
   Reinforcement learning

   Specify a fairness definition
- 3. See what guarantees we can make

# **Outline For This Talk**

 Specify a setting Reinforcement learning
 2. Specify a fairness definition

3. See what guarantees we can make

### Meritocratic fairness [JKMR16]:

Meritocratic fairness [JKMR16]: (whp), in every choice learning algorithm makes, for options a and b

Meritocratic fairness [JKMR16]: (whp), in every choice learning algorithm makes, for options a and b

# if quality(a) $\geq$ quality(b),

Meritocratic fairness [JKMR16]: (whp), in every choice learning algorithm makes, for options a and b

### if quality(a) $\geq$ quality(b), then P(choose a) $\geq$ P(choose b)








# In Our Example...



# In Our Example...



### if quality(a) $\geq$ quality(b), then P(choose a) $\geq$ P(choose b)

# if reward(a) $\geq$ reward(b),

### if quality(a) $\geq$ quality(b), then P(choose a) $\geq$ P(choose b)

# if reward(a) $\geq$ reward(b), then P(choose a) $\geq$ P(choose b)

### if quality(a) $\geq$ quality(b), then P(choose a) $\geq$ P(choose b)

# if reward(a) $\geq$ reward(b), then P(choose a) $\geq$ P(choose b)

Short-term

### if quality(a) $\geq$ quality(b), then P(choose a) $\geq$ P(choose b)

# if $Q^*(s, a) \ge Q^*(s, b)$ , then $P(s, a) \ge P(s, b)$

# if quality(a) $\geq$ quality(b), then P(choose a) $\geq$ P(choose b) Long-term if $Q^*(s, a) \ge Q^*(s, b)$ , then $P(s, a) \ge P(s, b)$

### Limitations

### Limitations

#### Minimal – more "necessary" than "sufficient"

### Limitations

#### Minimal – more "necessary" than "sufficient"

Assumes feedback reflects quality

#### Minimal – more "necessary" than "sufficient"

Minimal – more "necessary" than "sufficient"

Holds throughout learning process

Minimal – more "necessary" than "sufficient"

- Holds throughout learning process
- •Aligned with optimality!

### if quality(a) $\geq$ quality(b), then P(choose a) $\geq$ P(choose b)

#### Uniformly random exploration is fair

(since P(choose a) = P(choose b)) Uniformly random exploration is fair

#### Uniformly random exploration is fair

#### Uniformly random exploration is fair

## if quality(a) $\geq$ quality(b), then P(choose a) $\geq$ P(choose b)

#### Optimal exploitation is fair

Uniformly random exploration is fair

### if quality(a) $\geq$ quality(b), then P(choose a) $\geq$ P(choose b)

Optimal exploitation is fair (since quality(best)  $\geq$  quality(rest))

#### Uniformly random exploration is fair

#### Optimal exploitation is fair

#### Uniformly random exploration is fair

#### Need fair path between

#### Optimal exploitation is fair

# **Outline For This Talk**

- Specify a setting Reinforcement learning
   2. Specify a fairness definition Meritocratic fairness
- 3. See what guarantees we can make

# **Outline For This Talk**

 Specify a setting Reinforcement learning
 Specify a fairness definition Meritocratic fairness
 See what guarantees we can make

# **Our Performance Metric**

#### 



### Without Fairness

# Without Fairness

# Near-optimality takes poly(MDP parameters) steps

# Without Fairness

# Near-optimality takes poly(MDP parameters) steps

What does fairness cost?

### Lower Bound

### Lower Bound

Theorem: No fair algorithm can guarantee near-optimality in under exponential(# states) steps.

# Lower Bound Sketch
if  $Q^*(s, a) \ge Q^*(s, b)$ , then  $P(s, a) \ge P(s, b)$ 

if  $Q^*(s, a) \ge Q^*(s, b)$ , then  $P(s, a) \ge P(s, b)$ 

Fairness  $\rightarrow$  must explore randomly to learn Q<sup>\*</sup> values...

if  $Q^*(s, a) \ge Q^*(s, b)$ , then  $P(s, a) \ge P(s, b)$ 

Fairness  $\rightarrow$  must explore randomly to learn Q<sup>\*</sup> values...

...but sometimes random exploration does poorly



#### "Combination lock" MDP



#### "Combination lock" MDP



#### Exponential in # states

# How to get around this?

# How to get around this?

# How to get around this?

#### Idea: relax to approximate fairness

### if $Q^*(s, a) \ge Q^*(s, b)$ , then $\mathcal{L}(s, a) \ge \mathcal{L}(s, b)$

### if $Q^*(s, a) \ge Q^*(s, b)$ , then $\mathcal{L}(s, a) \ge \mathcal{L}(s, b)$

# if $Q^*(s, a) \ge Q^*(s, b)$ , then $\mathcal{L}(s, a) \geq \mathcal{L}(s, b)$ if $Q^*(s, a) + \alpha \ge Q^*(s, b)$ , then $\mathcal{L}(s, a) \geq \mathcal{L}(s, b)$

### if $Q^*(s, a) \ge Q^*(s, b)$ , then $\mathcal{L}(s, a) \geq \mathcal{L}(s, b)$ Approximate " "action" fairness if $Q^*(s, a) + \alpha \ge Q^*(s, b)$ , then $\mathcal{L}(s, a) \geq \mathcal{L}(s, b)$



### **Better!**

#### No longer exponential in *#* states



### **Better!**

#### No longer exponential in *#* states



### (Still) exponential in $1/(1-\gamma)$

# An Algorithm: Fair-E<sup>3</sup>

# An Algorithm: Fair-E<sup>3</sup>

### Start from E<sup>3</sup> [KS98]

# An Algorithm: Fair-E<sup>3</sup>

### Start from E<sup>3</sup> [KS98]

# Adapt to satisfy approximate-action fairness

### Organize world into "known" and "unknown" states

### Organize world into "known" and "unknown" states

"Known": good estimates of transitions, rewards, Q\* values . . .

### In unknown state: take random walk → state more known

#### In known state:

#### In known state:

# Either *fairly* exploit in known states for good reward

#### In known state:

# Either *fairly* exploit in known states for good reward

Or fairly explore to unknown quickly

# Approximate fairness makes everything trickier

# Approximate fairness makes everything trickier

### "Known" must be stronger

# Approximate fairness makes everything trickier

### "Known" must be stronger

Computing fair policies more delicate

# **Upper Bound**

# **Upper Bound**

Theorem: Fair-E<sup>3</sup> is approximate action fair and near-optimal in poly(all MDP parameters but  $\gamma$ ), exp(1/1- $\gamma$ ) steps

### **Collected Results**

## **Collected Results**

# Without fairness: near-optimality in poly(MDP parameters)
#### **Collected Results**

# Without fairness: near-optimality in poly(MDP parameters)

With fairness: exp(# states)

#### **Collected Results**

# Without fairness: near-optimality in poly(MDP parameters)

With fairness: exp(# states)

With approximate fairness: exp(discount factor)

• Fair ML matters!

- Fair ML matters!
- Proposed meritocratic fairness, studied in RL setting

- Fair ML matters!
- Proposed meritocratic fairness, studied in RL setting
- Proved separations between "unfair", fair, and approximately fair RL

- Fair ML matters!
- Proposed meritocratic fairness, studied in RL setting
- Proved separations between "unfair", fair, and approximately fair RL
  Can we do better?

- Fair ML matters!
- Proposed meritocratic fairness, studied in RL setting
- Proved separations between "unfair", fair, and approximately fair RL
  Can we do better? Thanks!

Paper: arxiv.org/abs/1611.03071

#### References

[JKMR16] "Fairness in Learning: Classic and Contextual Bandits" Matthew Joseph, Michael Kearns, Jamie Morgenstern, and Aaron Roth. NIPS 2016

[KSO2] "Near-Optimal Reinforcement Learning in Polynomial Time" Michael Kearns and Satinder Singh. ICML 1998