

Real World Reinforcement Learning



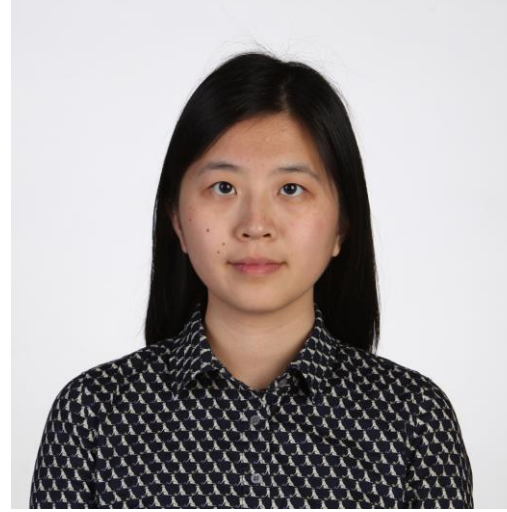
Jack
Gerrits



Rodrigo
Kumpera



John
Langford



Cheng
Tan



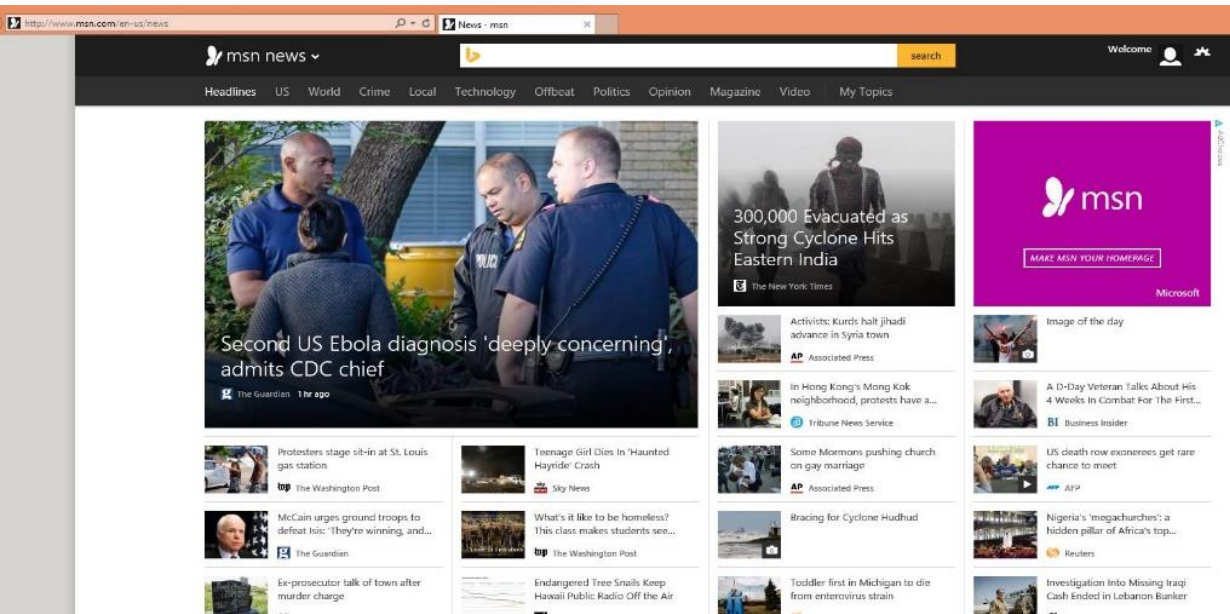
Alexey
Taymanov

ICML RWRL Workshop, June 9, 2019

Slides/references at <https://vowpalwabbit.github.io/icml2019/>

Which News?

Why?



28% lift

Which Game?



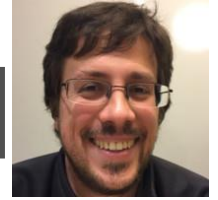
40% lift

Outline of the afternoon

1. Core ideas [John Langford]



2. The Personalizer System [Rodrigo Kumpera]

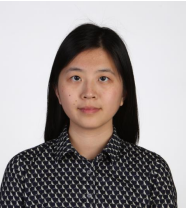


Break

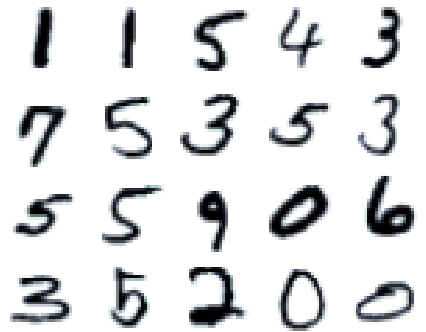
1. Hands-on using it [Alexey Taymanov]



2. Hands-on counterfactual Evaluation [Cheng Tan]



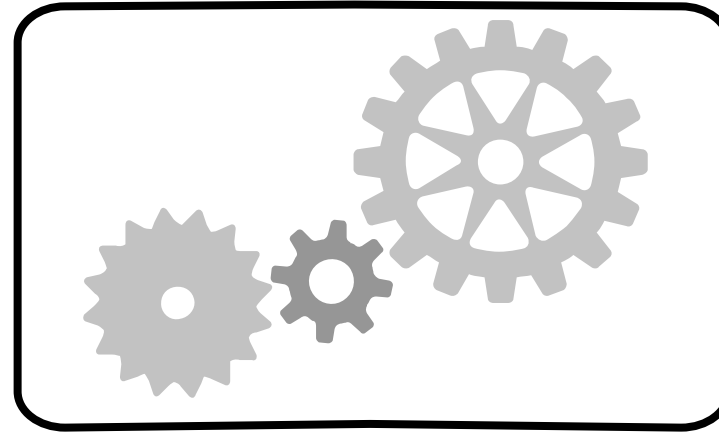
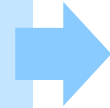
The Baseline -- Supervised Learning



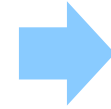
Training examples

1	1	5	4	3
7	5	3	5	3
5	5	9	0	6
3	5	2	0	0

Training labels



Supervised Learner



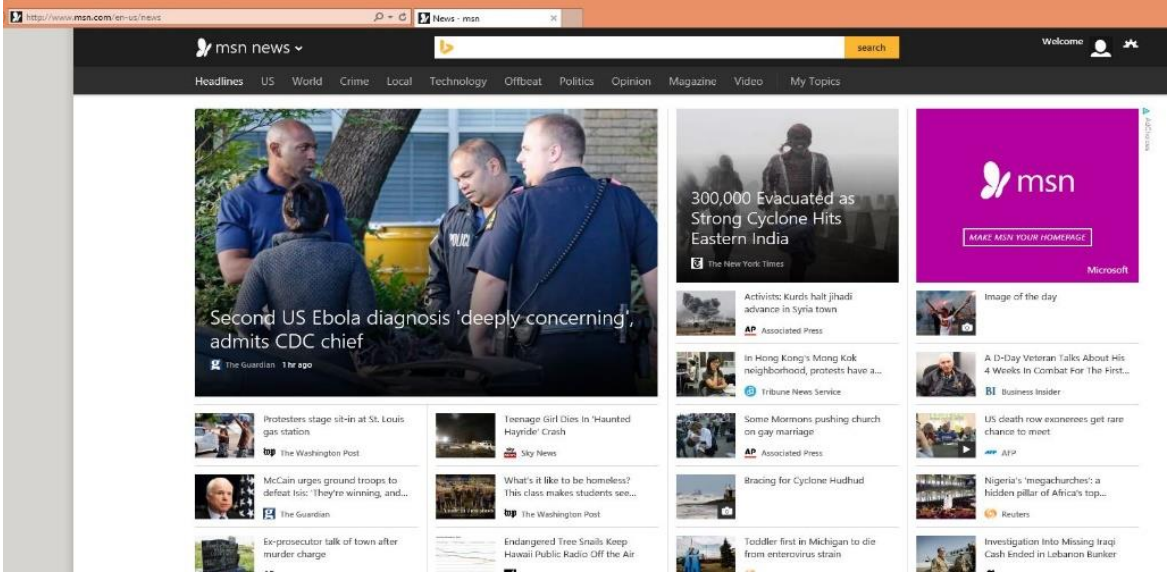
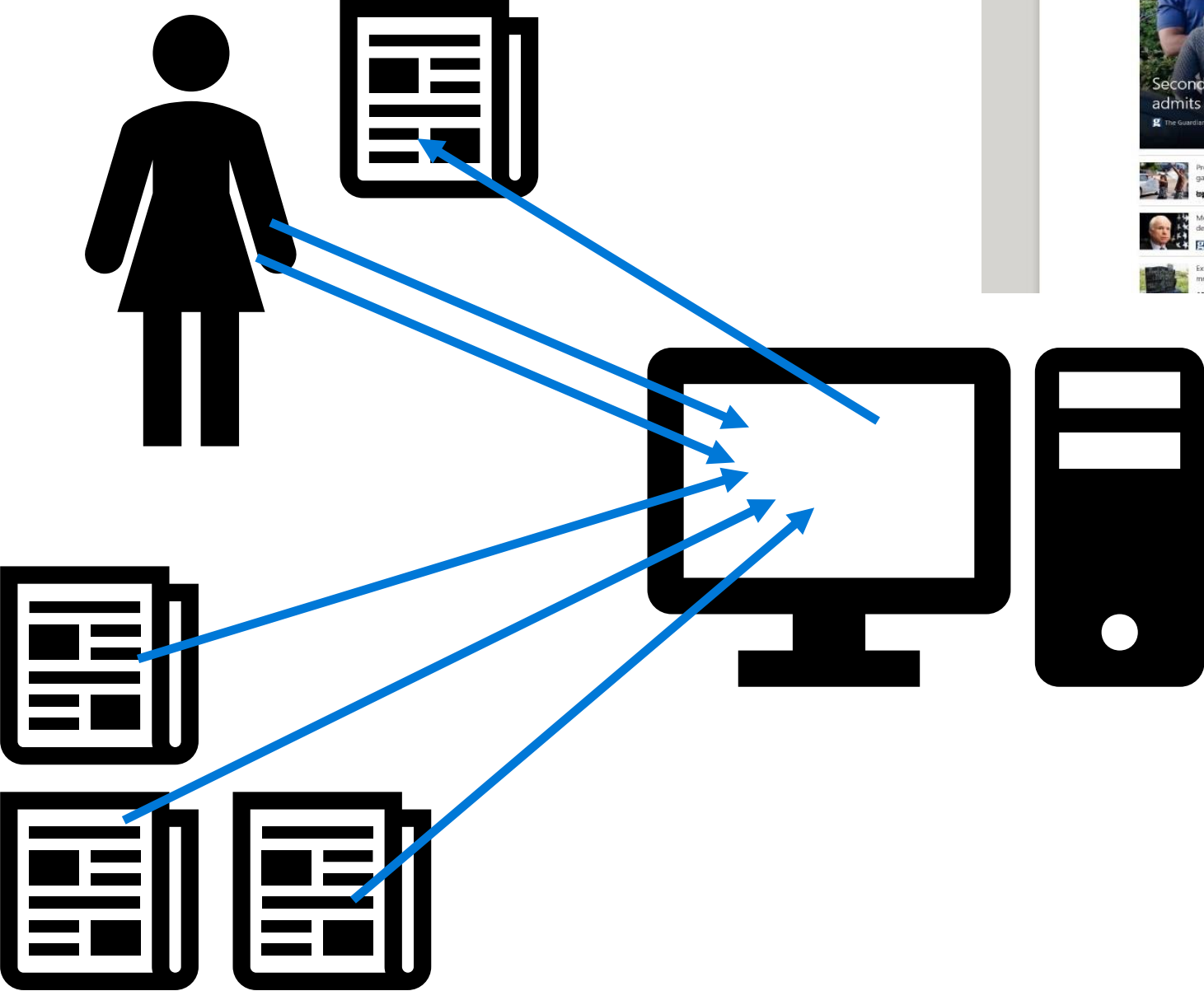
Accurate digit
classifier

2

Supervised Learning is cool



How about news?



A standard pipeline

1. Collect $(user, article)$ information.
2. Build $features(user, article)$
3. Hire editor to judge $relevance(user, article)$
4. Learn $\widehat{rel}(features(user, article))$
5. Act: $\arg \max_{\{articles\}} \widehat{rel}(features(user, article))$
6. Deploy in A/B test for 2 weeks
7. A/B test fails 😞

A standard pipeline

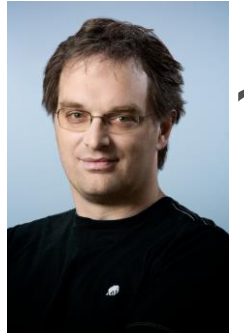
1. Collect $(user, article, click)$ information.
2. Build $features(user, article)$
3. Learn $\hat{P}(click|features(user, article))$
4. Act: $\arg \max_{\{articles\}} \hat{P}(click|features(user, article))$
5. Deploy in A/B test for 2 weeks
6. A/B test fails 😞 Why?

Q: What goes wrong?

Is Ukraine



interesting to John



?

A: Need Right Signal for Right Answer

What goes wrong?

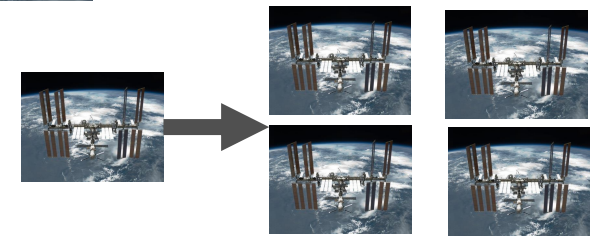
What is the probability of click on a food article



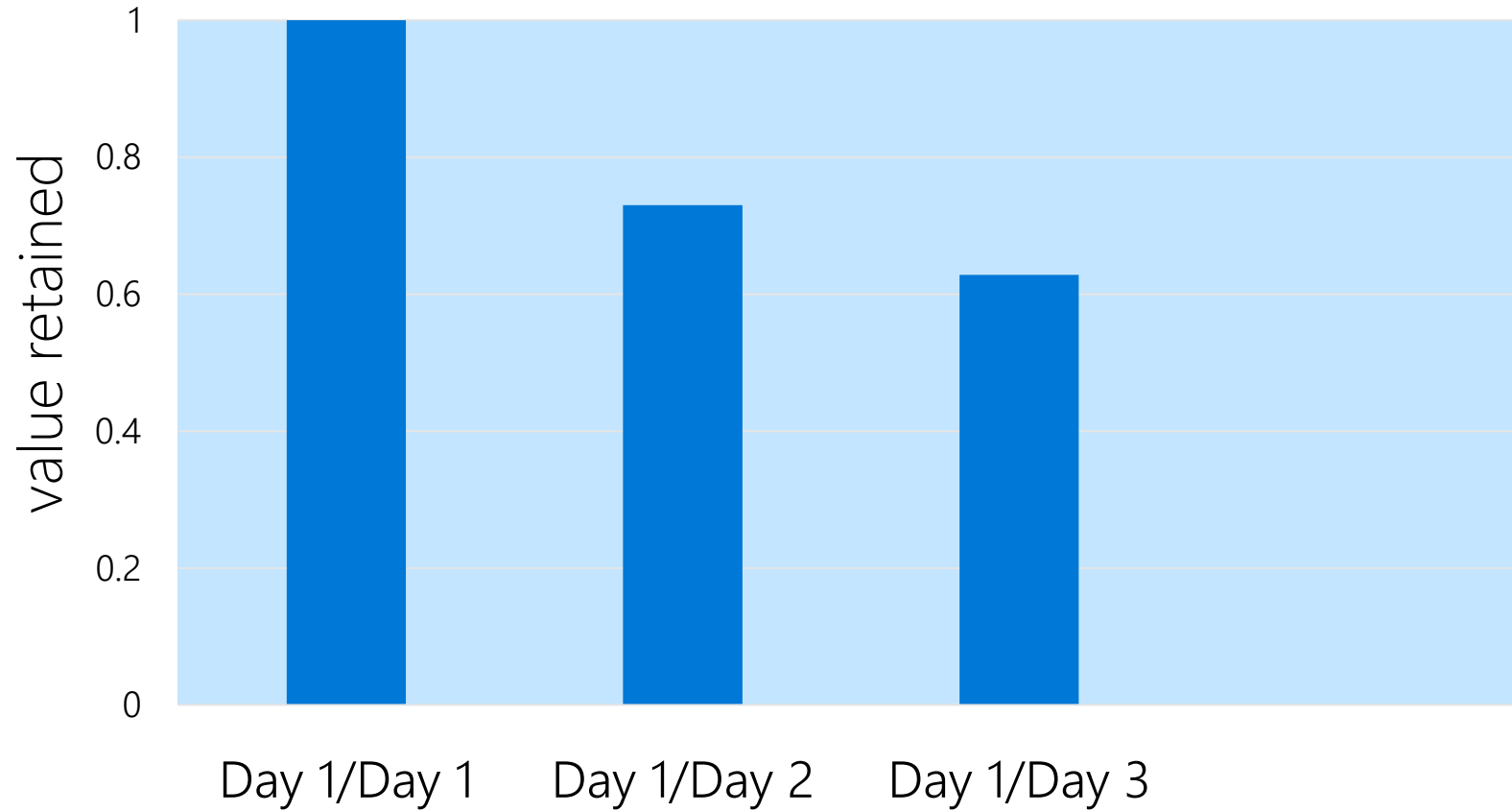
If you only display a space article?



We must avoid “self-fulfilling prophecy”



What else goes wrong?



The world changes!

Can we optimize for the best outcome?

Amongst a given set of choices?

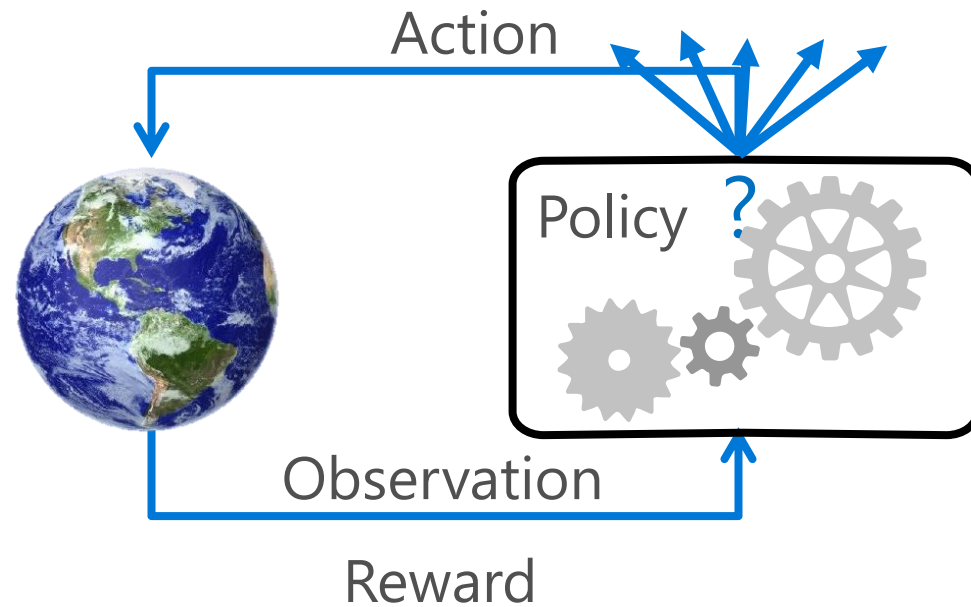
For what matters to individuals?

Without self-fulfilling prophecies?

With real-time learning?



Reinforcement Learning can do this!



Goal: Find a policy maximizing the sum of rewards

Q: One last Why...



~~AI: A function programmed with data~~

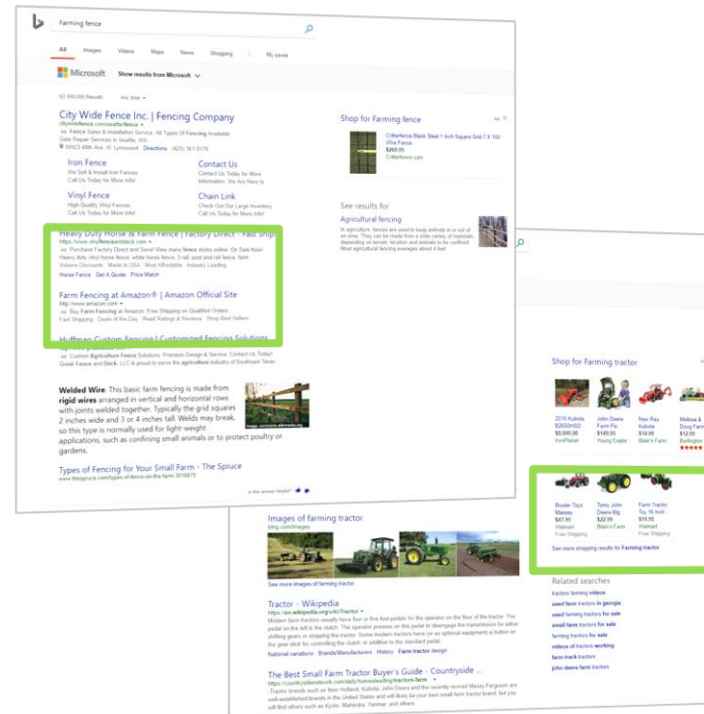
AI: An economically viable digital agent that explores, learns, and acts

Content

Layout

Creative

“Book Your Vacation to Hawaii”



Wellness

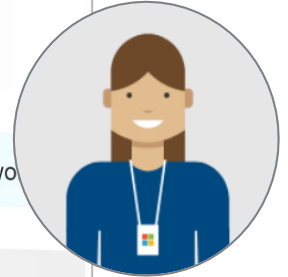


[ZKZ '09, SLLSPM '11, NSTWCSM '14,
PGCRRH '14, NHS '15, KHSBATM '15,
HFKMTY '16]

Bots

I'm Microsoft's Virtual Agent. I'd love to help you. You can also ask to talk to a person at any time. Please briefly describe your issue below.

My printer isn't wo



Check power and connection

If you are unable to print or connect to your printer in Windows 10, first try this:

1. Make sure that your printer is plugged into the power supply and turned on.
2. Check the USB connection (for wired printers) or the wireless connection (for wireless printers).

Did that solve the problem?

[Yes](#)

[No, show solution 2 of 5](#)

Other Real-world Applications

Ad Choice: [BPQCCPRSS '12]

Ad Format: [TRSA '13]

Education: [MLLBP '14]

Music Rec: [WWHW '14]

Robotics: [PG '16]

Formalism: Contextual Bandits

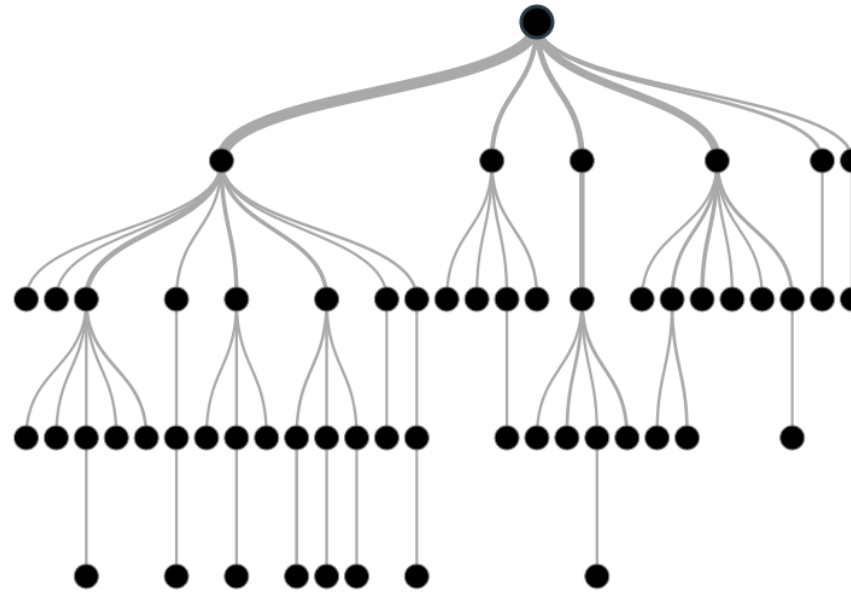
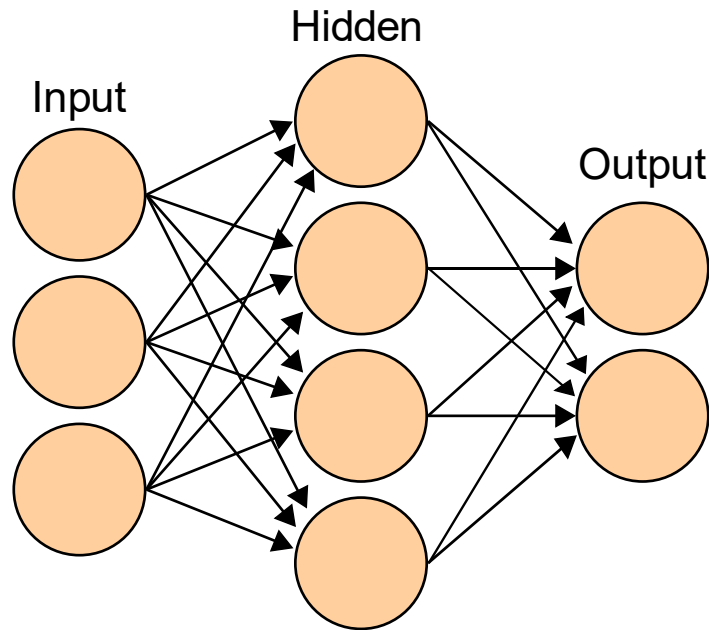
Repeatedly:

1. Observe features x
2. Choose action $a \in A$
3. Observe reward r

Goal: Maximize expected reward

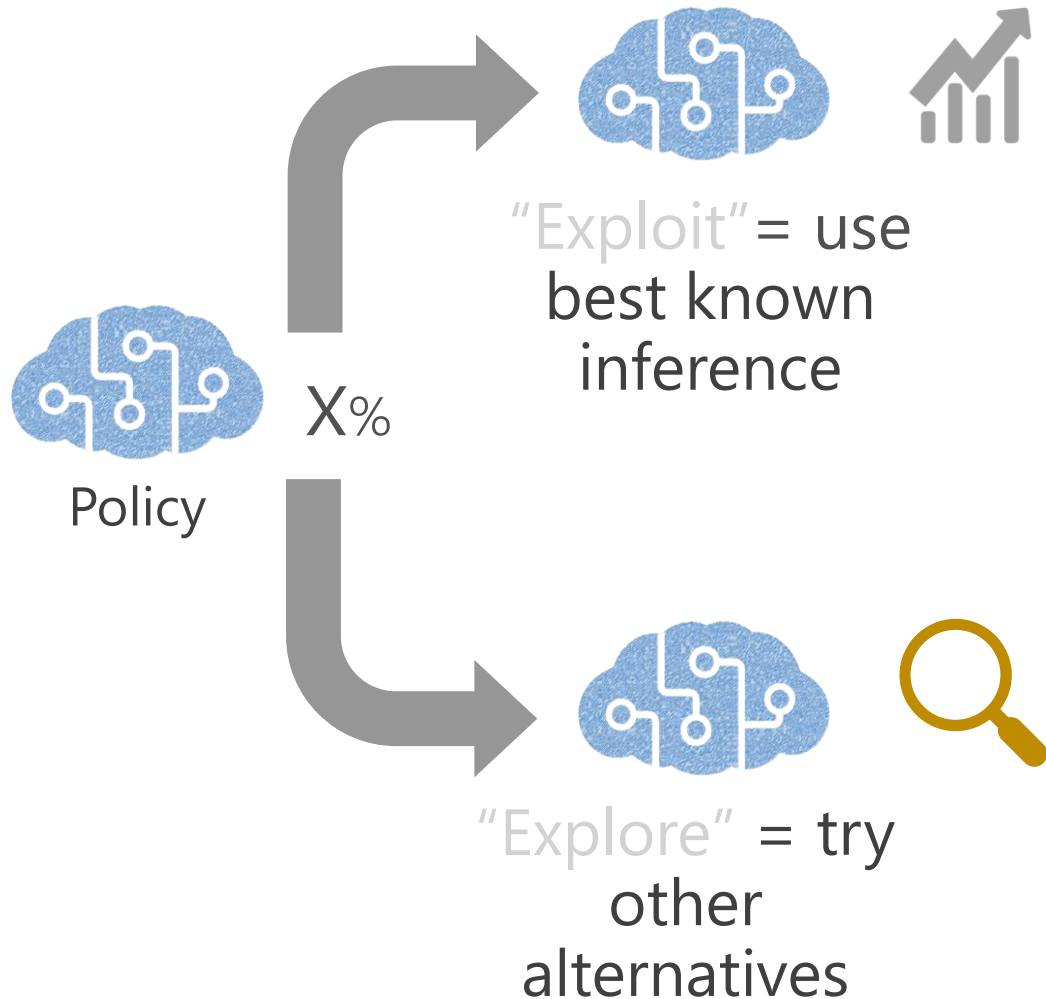
Policies

Policy maps features to actions.



Policy = Classifier that *acts*.

Why does it work?

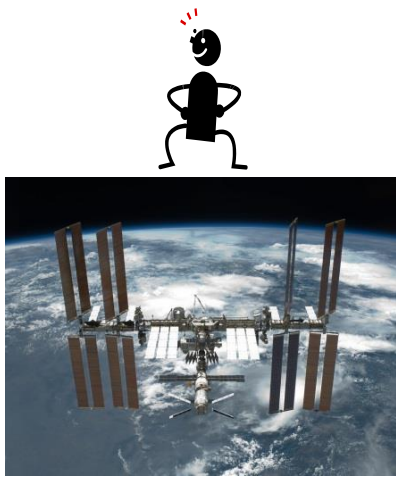


Exploit for performance

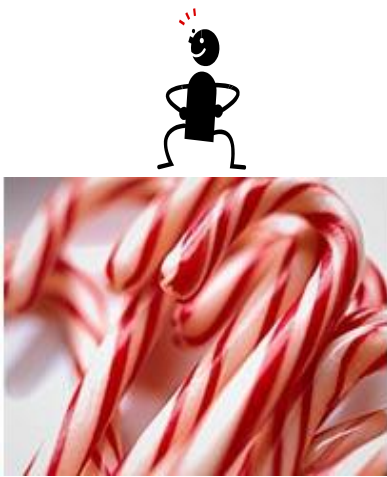
"How much should I explore to discover how to best perform?"

Explore to discover new things

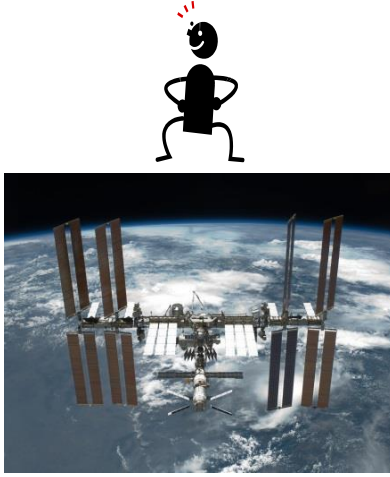
Counterfactual Evaluation



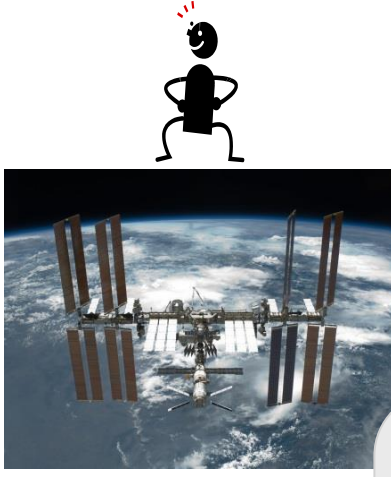
Read



Read



Ignored



Read

...

Tests can use the same events!

Later evaluate Location rule:

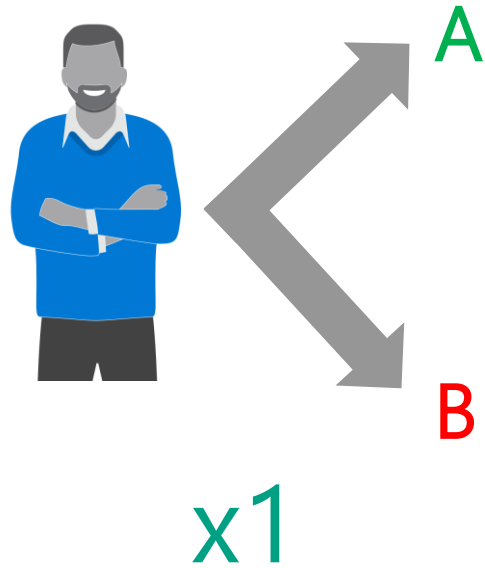
Teacher
Texas

Engineer
Seattle

Engineer
Seattle

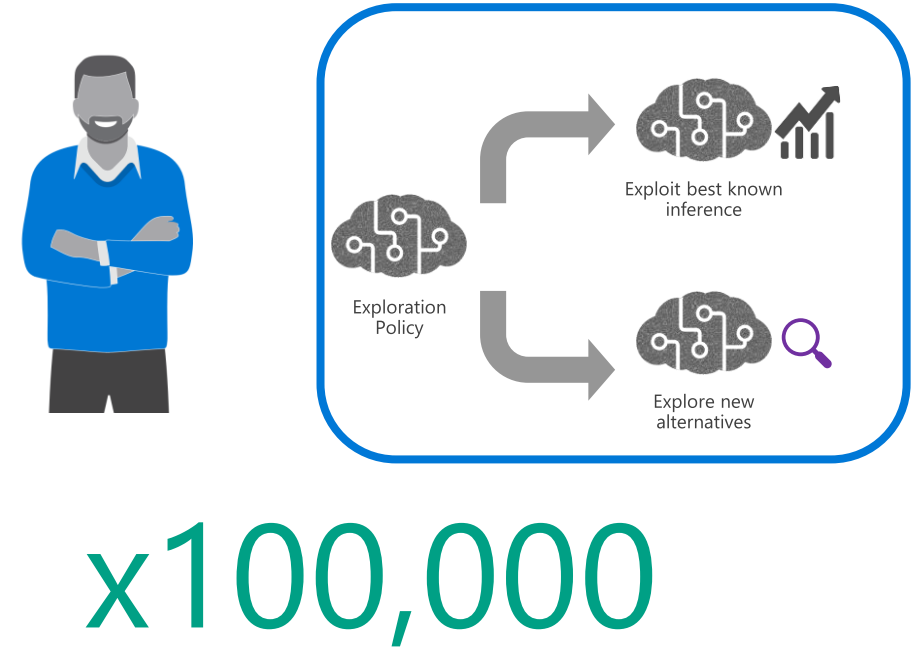
Engineer
Texas

A/B Testing vs. Counterfactual Evaluation



A/B Test:

1. Design the Right Experiment,
2. Test online once
3. Start over




Offline Experiment:

1. Use models that exploit and explore
2. Record User Interaction
3. Find the policy and model that fits reality

Inverse Propensity Score(IPS) [HT '52]

Given experience $\{(x, a, p, r)\}$ and a policy $\pi: x \rightarrow a$, how good is π ?

$$V_{IPS}(\pi) = \frac{1}{n} \sum_{(x,a,p,r)} \frac{r I(\pi(x) = a)}{p}$$


Propensity Score

What do we know about IPS?

Theorem: For all π , for all $D(x, \vec{r})$

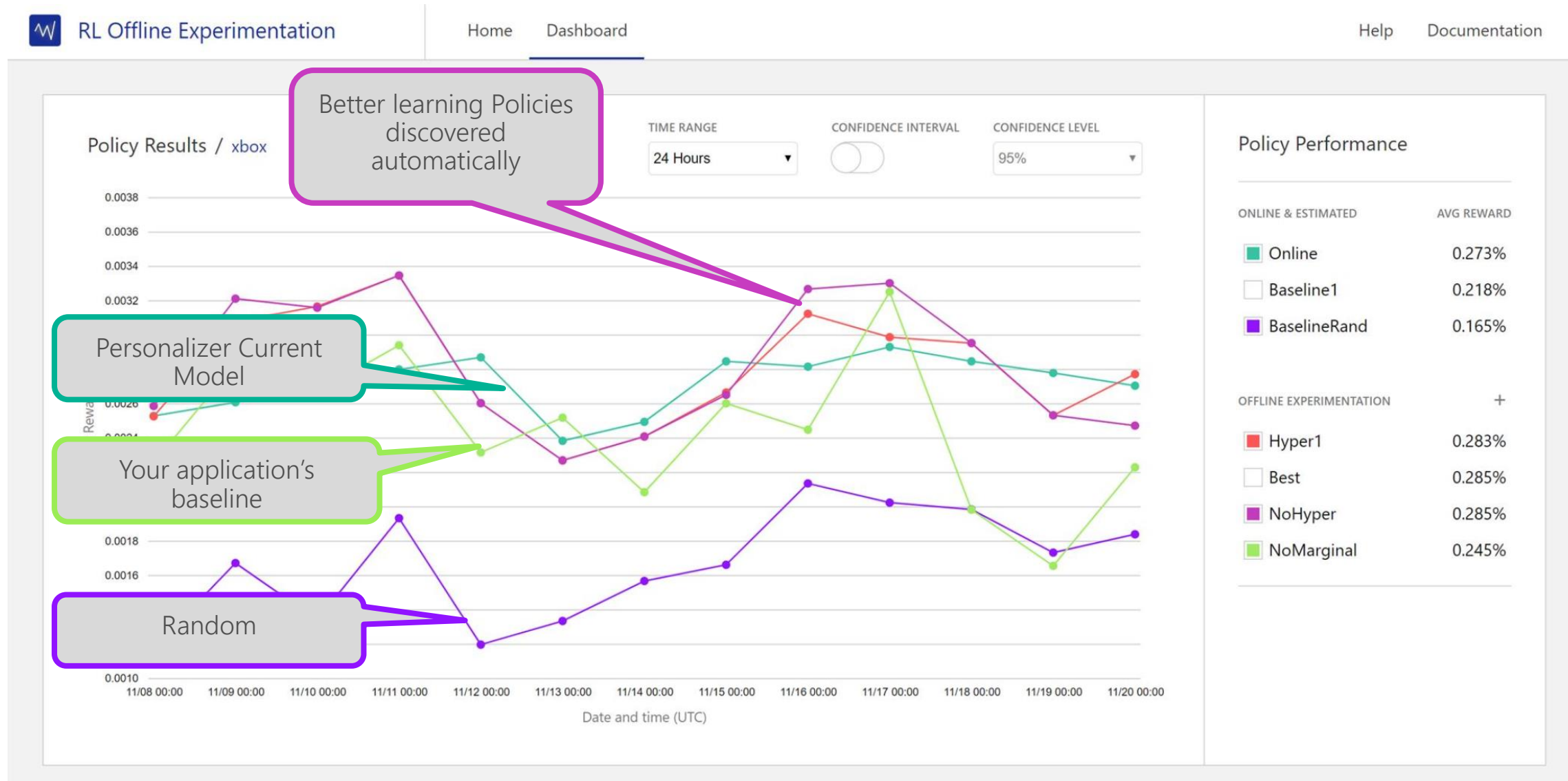
$$\mathbb{E} \left[r_{\pi(x)} \right] = \mathbb{E}[V_{\text{IPS}}(\pi)] = \mathbb{E} \left[\frac{1}{n} \sum_{(x,a,p,r)} \frac{r I(\pi(x)=a)}{p} \right]$$

Proof: For all (x, \vec{r}) , $E_{a \sim \vec{p}} \left[\frac{r_a I(\pi(x)=a)}{p_a} \right]$

$$= \sum_a p_a \frac{r_a I(\pi(x)=a)}{p_a}$$

$$= r_{\pi(x)}$$

Why Explore? You can do learning



Better Evaluation Techniques

Double Robust: [DLL '11]

Weighted IPS: [K '92, SJ '15]

Clipping: [BL '08]

Empirical Likelihood: [MKL '19]

Learning from Exploration [Z 03]

Given Data $\{(x, a, p, r)\}$ how to maximize $E[r_{\pi(x)}]$?

Maximize $E[V_{IPS}(\pi)]$ instead!

$$r_a = \begin{cases} r/p & \text{if } \pi(x) = a \\ 0 & \text{otherwise} \end{cases}$$

Equivalent to:

$$r'_a = \begin{cases} 1 & \text{if } \pi(x) = a \\ 0 & \text{otherwise} \end{cases}$$

with importance weight $\frac{r}{p}$

Importance weighted multiclass classification!

Better Learning from Exploration Data

Policy Gradient: [W '92]

Offset Tree: [BL '09]

Double Robust for learning: [DLL '11]

Multitask Regression: [BAL '18]

Weighted IPS for learning: [SJ '15]

Evaluating Online Learning

Problem: How do you evaluate an online learning algorithm
Offline?

Answer: Use Progressive Validation [BKL '99, CCG '04]



Theorem:

- 1) Expected PV value = Uniform expected policy value.
- 2) Trust like a **test set error**.

How do you do Exploration?

Simplest Algorithm: ϵ -greedy.

With probability ϵ act uniform random

With probability $1 - \epsilon$ act greedily

Better Exploration Algorithms

Better algorithms maintain ensemble and explore amongst actions of this ensemble.

Thompson Sampling: [T '33]

EXP4: [ACFS '02]

Epoch Greedy: [LZ '07]

Polytime: [DHKKLRZ '11]

Cover&Bag: [AHKLLS '14]

Bootstrap: [EK '14]

Evaluating Exploration Algorithms

Problem: How do you take the choice of examples acquired by an exploration algorithm into account?

Answer: Rejection Sample from history. [DELL '12]



Theorem: Realized history is unbiased up to length observed.

Better versions: [DELL '14]

More Research Details!

ICML tutorial: <http://hunch.net/~rwil>

John's Spring 2017 Cornell Tech class <http://hunch.net/~mltf>

Alekh's 2019 class http://alekhagarwal.net/bandits_and_rl/

Take-aways

- 1) Good fit for many problems
- 2) Fundamental questions have useful answers