

# **Exploration in RL: Contextual Bandit**

# Recap: MAB

## Interactive learning process:

For  $t = 0 \rightarrow T - 1$

(# based on historical information)

1. Learner pulls arm  $I_t \in \{1, \dots, K\}$
2. Learner observes an i.i.d reward  $r_t \sim \nu_{I_t}$  of arm  $I_t$

# Recap: MAB

## Interactive learning process:

For  $t = 0 \rightarrow T - 1$

(# based on historical information)

1. Learner pulls arm  $I_t \in \{1, \dots, K\}$
2. Learner observes an i.i.d reward  $r_t \sim \nu_{I_t}$  of arm  $I_t$

## Learning metric:

$$\text{Regret}_T = T\mu^\star - \sum_{t=0}^{T-1} \mu_{I_t}$$

# Recap: MAB

## Interactive learning process:

For  $t = 0 \rightarrow T - 1$

(# based on historical information)

1. Learner pulls arm  $I_t \in \{1, \dots, K\}$
2. Learner observes an i.i.d reward  $r_t \sim \nu_{I_t}$  of arm  $I_t$

## Learning metric:

$$\text{Regret}_T = T\mu^* - \sum_{t=0}^{T-1} \mu_{I_t}$$

Arm distributions are fixed across learning..

Question for Today:

Incorporate contexts into the interactive learning framework

## Outline for today:

1. Introduction of the model

2. Algorithm

3. Theory and some practical considerations

# Make the framework Context Dependent:

**Interactive learning process:**

For  $t = 0 \rightarrow T - 1$

**1. A new context  $x_t \in \mathcal{X}$  appears**

$\subseteq \mathbb{R}^d$

# Make the framework Context Dependent:

## Interactive learning process:

For  $t = 0 \rightarrow T - 1$

**1. A new context  $x_t \in \mathcal{X}$  appears**

(# based on context  $x_t$  and  
historical information)

2. Learner picks action  $a_t \in \mathcal{A}$



# Make the framework Context Dependent:

## Interactive learning process:

For  $t = 0 \rightarrow T - 1$

1. A new context  $x_t \in \mathcal{X}$  appears ✓

(# based on context  $x_t$  and historical information)

2. Learner picks action  $a_t \in \mathcal{A}$

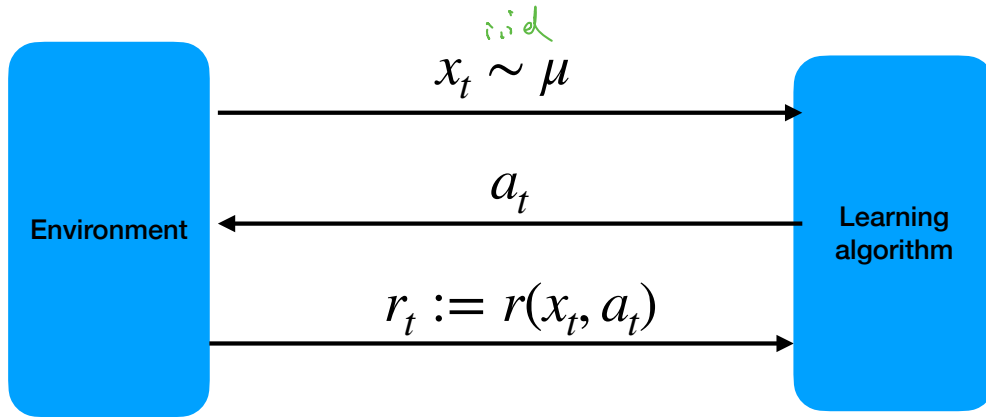
3. Learner observes an reward  $r_t := r(x_t, a_t)$

Deterministic

Reward is context and arm dependent now!

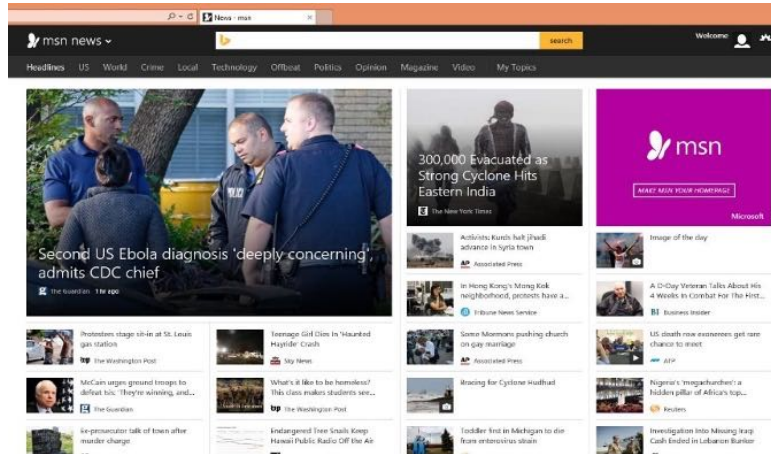
# Make the framework Context Dependent:

**Interactive learning process:**



# Examples:

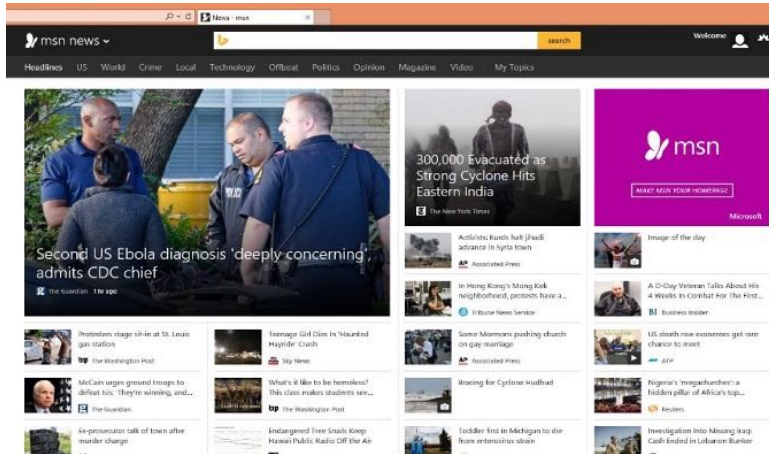
## Personalize recommendation system



# Examples:

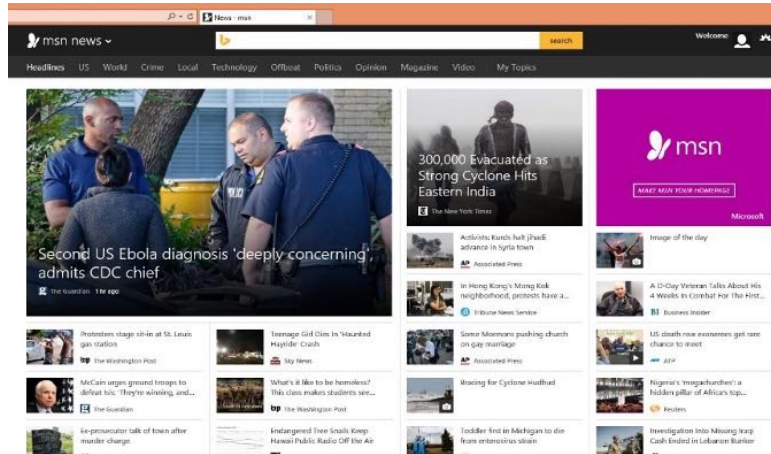
## Personalize recommendation system

**Context:** user's information (e.g., history, health conditions, age, height, weight, job type, etc)



# Examples:

## Personalize recommendation system

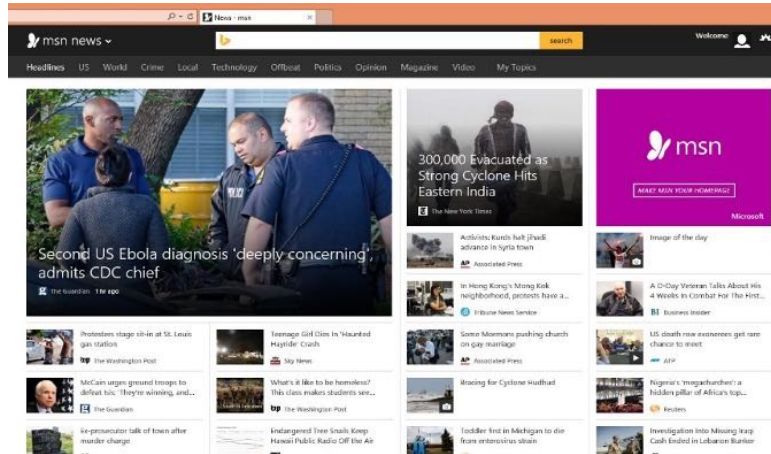


**Context:** user's information (e.g., history health conditions, age, height, weight, job type, etc)

**Decisions (arms):** news articles

# Examples:

## Personalize recommendation system



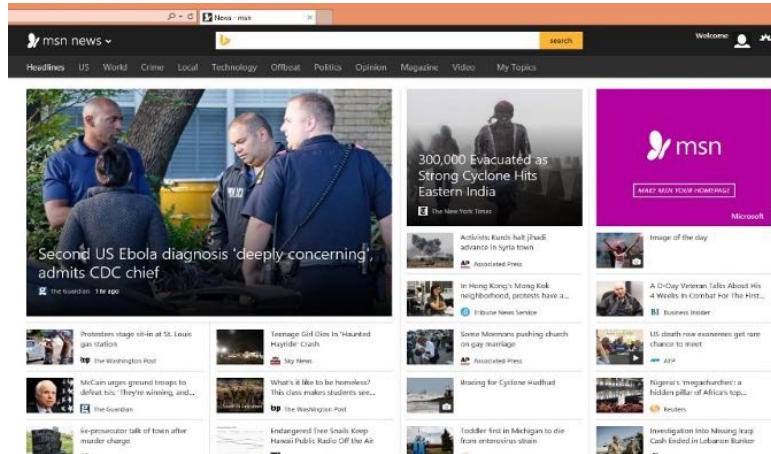
**Context:** user's information (e.g., history health conditions, age, height, weight, job type, etc)

**Decisions (arms):** news articles

**Goal:** learn to maximize user click rate

# Examples:

## Personalize recommendation system



**Context:** user's information (e.g., history health conditions, age, height, weight, job type, etc)

**Decisions (arms):** news articles

**Goal:** learn to maximize user click rate

Different users have different preferences on news, so need to personalize

Equivalently, it is an MDP with  $H = 1$

Finite horizon MDP with  $H = 1$

$$\mathcal{M} = \{\mathcal{X}, \mathcal{A}, r, H = 1, \mu\}$$

*Handwritten annotations:*  
A green arrow points from the word "state" to the  $\mathcal{X}$  in the set definition.  
A green double-headed arrow is positioned below the  $\mu$  in the set definition.



Equivalently, it is an MDP with  $H = 1$

Finite horizon MDP with  $H = 1$

$$\mathcal{M} = \{\mathcal{X}, \mathcal{A}, r, H = 1, \mu\}$$

Objective function:

$$\max_{\pi \in \Pi} \mathbb{E}_{x \sim \mu} [r(x, \pi(x))]$$

Equivalently, it is an MDP with  $H = 1$

Finite horizon MDP with  $H = 1$

$$\mathcal{M} = \{\mathcal{X}, \mathcal{A}, r, H = 1, \mu\}$$

Objective function:

$$\max_{\pi \in \Pi} \mathbb{E}_{x \sim \mu} \left[ r(x, \pi(x)) \right]$$

For simplicity, we assume reward is deterministic;  
The challenge is really from randomness in contexts

# The Regret Metric

Fix a policy class  $\Pi$  ( think about  $\pi$  as a classifier from  $x \rightarrow a$ )

Denote optimal policy  $\pi^* = \arg \max_{\pi \in \Pi} \mathbb{E}_{x \sim \mu} r(x, \pi(x))$

$\psi$   $\mathcal{A}$

# The Regret Metric

Fix a policy class  $\Pi$  ( think about  $\pi$  as a classifier from  $x \rightarrow a$ )

Denote optimal policy  $\pi^* = \arg \max_{\pi \in \Pi} \mathbb{E}_{x \sim \mu} r(x, \pi(x))$

Every iteration, learner has a policy  $\pi^t \in \Pi$   
(Recommends  $a_t = \pi^t(x_t)$ , receives reward  $r_t := r(x_t, \pi(x_t))$ )

# The Regret Metric

Fix a policy class  $\Pi$  ( think about  $\pi$  as a classifier from  $x \rightarrow a$ )

Denote optimal policy  $\pi^\star = \arg \max_{\pi \in \Pi} \mathbb{E}_{x \sim \mu} r(x, \pi(x))$

Every iteration, learner has a policy  $\pi^t \in \Pi$

(Recommends  $a_t = \pi^t(x_t)$ , receives reward  $r_t := r(x_t, \pi(x_t))$ )

$$\text{Regret}_T = T \mathbb{E}_{x \sim \mu} [r(x, \pi^\star(x))] - \sum_{t=0}^{T-1} \mathbb{E}_{x \sim \mu} [r(x, \pi^t(x))]$$

# The Regret Metric

Fix a policy class  $\Pi$  ( think about  $\pi$  as a classifier from  $x \rightarrow a$ )

Denote optimal policy  $\pi^{\star} = \arg \max_{\pi \in \Pi} \mathbb{E}_{x \sim \mu} r(x, \pi(x))$

Every iteration, learner has a policy  $\pi^t \in \Pi$

(Recommends  $a_t = \pi^t(x_t)$ , receives reward  $r_t := r(x_t, \pi(x_t))$ )

$$\text{Regret}_T = T \mathbb{E}_{x \sim \mu} [r(x, \pi^{\star}(x))] - \sum_{t=0}^{T-1} \mathbb{E}_{x \sim \mu} [r(x, \pi^t(x))]$$

Total expected reward if we always  
uses  $\pi^{\star}$  to recommend

# The Regret Metric

Fix a policy class  $\Pi$  ( think about  $\pi$  as a classifier from  $x \rightarrow a$ )

Denote optimal policy  $\pi^* = \arg \max_{\pi \in \Pi} \mathbb{E}_{x \sim \mu} r(x, \pi(x))$

Every iteration, learner has a policy  $\pi^t \in \Pi$

(Recommends  $a_t = \pi^t(x_t)$ , receives reward  $r_t := r(x_t, \pi(x_t))$ )

$$\text{Regret}_T = T \mathbb{E}_{x \sim \mu} [r(x, \pi^*(x))] - \sum_{t=0}^{T-1} \mathbb{E}_{x \sim \mu} [r(x, \pi^t(x))]$$

Total expected reward if we always  
uses  $\pi^*$  to recommend

Total expected reward of our learned  
sequence of policies

## Outline for today:

 1. Introduction of the model

2. Algorithm

3. Theory and some practical considerations



# Ingredient 1: Importance Weighting

The key challenging here is that we observe  $r_t := r(x_t, a_t)$ ,  
but we do not know  $r(x_t, a)$  for  $a \neq a_t$

# Ingredient 1: Importance Weighting

The key challenging here is that we observe  $r_t := r(x_t, a_t)$ ,  
but we do not know  $r(x_t, a)$  for  $a \neq a_t$

Importance weighting actually allows us to get  
**unbiased estimate for ALL actions!**

# Ingredient 1: Importance Weighting

The key challenging here is that we observe  $r_t := r(x_t, a_t)$ ,  
but we do not know  $r(x_t, a)$  for  $a \neq a_t$

Importance weighting actually allows us to get  
**unbiased estimate for ALL actions!**

Assume  $a_t \sim p$  ( $p \in \Delta(\mathcal{A})$ ), and we log  $p(a_t)$ , receive  $r_t = r(x_t, a_t)$ ,

# Ingredient 1: Importance Weighting

The key challenging here is that we observe  $r_t := r(x_t, a_t)$ ,  
but we do not know  $r(x_t, a)$  for  $a \neq a_t$

Importance weighting actually allows us to get  
**unbiased estimate for ALL actions!**

Assume  $a_t \sim p$  ( $p \in \Delta(\mathcal{A})$ ), and we log  $p(a_t)$ , receive  $r_t = r(x_t, a_t)$ ,

For all  $a \in \mathcal{A}$ , define  $\hat{\mathbf{r}}[a] = \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)}$ ,

# Ingredient 1: Importance Weighting

The key challenging here is that we observe  $r_t := r(x_t, a_t)$ ,  
but we do not know  $r(x_t, a)$  for  $a \neq a_t$

Importance weighting actually allows us to get  
**unbiased estimate for ALL actions!**

Assume  $a_t \sim p$  ( $p \in \Delta(\mathcal{A})$ ), and we log  $p(a_t)$ , receive  $r_t = r(x_t, a_t)$ ,

For all  $a \in \mathcal{A}$ , define  $\hat{\mathbf{r}}[a] = \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)}$ ,

$$\hat{\mathbf{r}} := \begin{bmatrix} 0 \\ 0 \\ \vdots \\ r_t/p(a_t) \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow a_t$$

# Ingredient 1: Importance Weighting

The key challenging here is that we observe  $r_t := r(x_t, a_t)$ ,  
but we do not know  $r(x_t, a)$  for  $a \neq a_t$

Importance weighting actually allows us to get  
**unbiased estimate for ALL actions!**

Assume  $a_t \sim p$  ( $p \in \Delta(\mathcal{A})$ ), and we log  $p(a_t)$ , receive  $r_t = r(x_t, a_t)$ ,

For all  $a \in \mathcal{A}$ , define  $\hat{\mathbf{r}}[a] = \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)}$ ,

$$\mathbb{E}_{a_t \sim p} \hat{\mathbf{r}}[a] = r(x_t, a), \forall a \in \mathcal{A}$$

$$\hat{\mathbf{r}} := \begin{bmatrix} 0 \\ 0 \\ \dots \\ r_t/p(a_t) \\ 0, \\ \dots \\ 0 \end{bmatrix}$$

$\mathbb{E}[\hat{\mathbf{r}}] = r(x_t, \cdot)$

# Proving Importance Weighting

Assume  $a_t \sim p$  ( $p \in \Delta(\mathcal{A})$ ), and we log  $p_t(a)$ , receive  $r_t = r(x_t, a_t)$ ,

For all  $a \in \mathcal{A}$ , define  $\hat{r}[a] = \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)}$ , we have:  $\mathbb{E}_{a_t \sim p} \hat{r}[a] = r(x_t, a), \forall a \in \mathcal{A}$

# Proving Importance Weighting

Assume  $a_t \sim p$  ( $p \in \Delta(\mathcal{A})$ ), and we log  $p_t(a)$ , receive  $r_t = r(x_t, a_t)$ ,

For all  $a \in \mathcal{A}$ , define  $\hat{r}[a] = \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)}$ , we have:  $\mathbb{E}_{a_t \sim p} \hat{r}[a] = r(x_t, a), \forall a \in \mathcal{A}$

*Randomness  
is from  $a_t$*

Consider any  $a \in \mathcal{A}$  :



# Proving Importance Weighting

Assume  $a_t \sim p$  ( $p \in \Delta(\mathcal{A})$ ), and we log  $p_t(a)$ , receive  $r_t = r(x_t, a_t)$ ,

For all  $a \in \mathcal{A}$ , define  $\hat{r}[a] = \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)}$ , we have:  $\mathbb{E}_{a_t \sim p} \hat{r}[a] = r(x_t, a), \forall a \in \mathcal{A}$

Consider any  $a \in \mathcal{A}$  :  $\mathbb{E}_{a_t \sim p} [\hat{r}[a]] = r(x_t, a)$

$$\mathbb{E}_{a_t \sim p} \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)} = \sum_{a_t \in \mathcal{A}} p(a_t) \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)} = \mathbb{E}_{a_t \sim p}$$

# Proving Importance Weighting

Assume  $a_t \sim p$  ( $p \in \Delta(\mathcal{A})$ ), and we log  $p_t(a)$ , receive  $r_t = r(x_t, a_t)$ ,

For all  $a \in \mathcal{A}$ , define  $\hat{r}[a] = \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)}$ , we have:  $\mathbb{E}_{a_t \sim p} \hat{r}[a] = r(x_t, a), \forall a \in \mathcal{A}$

Consider any  $a \in \mathcal{A}$ :

$$\begin{aligned}\mathbb{E}_{a_t \sim p} \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)} &= \sum_{a_t \in \mathcal{A}} p(a_t) \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)} \\ &= \cancel{p(a)} \frac{r(x_t, a)}{\cancel{p(a)}} = \underline{r(x_t, a)}\end{aligned}$$

← fires when  $a_t = a$

$$\mathbb{E}_{a_t \sim p} [\hat{r}[a]] = r(x_t, a)$$

# Ingredient 2: Reward-sensitive Classification Oracle

Recall classic classification:

# Ingredient 2: Reward-sensitive Classification Oracle

Recall classic classification:

$$\{x_i, y_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, y_i \in \{1, \dots, |\mathcal{A}|\}$$

$\uparrow$                        $\uparrow$   
feature                      label

# Ingredient 2: Reward-sensitive Classification Oracle

Recall classic classification:

$$\{x_i, y_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, y_i \in \{1, \dots, |\mathcal{A}|\}$$

$$\arg \max_{\pi \in \Pi} \sum_{i=1}^N \mathbf{1}\{\pi(x_i) = y_i\}$$

$$\begin{aligned} & \max_{\pi} \sum_{i=1}^N \mathbf{1}\{\pi(x_i) = y_i\} \\ \Leftrightarrow & \min_{\pi} \sum_{i=1}^N \mathbf{1}\{\pi(x_i) \neq y_i\} \end{aligned}$$

# Ingredient 2: Reward-sensitive Classification Oracle

Recall classic classification:

$$\{x_i, y_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, y_i \in \{1, \dots, |\mathcal{A}|\}$$

$$\arg \max_{\pi \in \Pi} \sum_{i=1}^N \mathbf{1}\{\pi(x_i) = y_i\}$$

**Let's generalize it to Reward-Sensitive Classification:**

# Ingredient 2: Reward-sensitive Classification Oracle

Recall classic classification:

$$\{x_i, y_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, y_i \in \{1, \dots, |\mathcal{A}|\}$$

$$\arg \max_{\pi \in \Pi} \sum_{i=1}^N \mathbf{1}\{\pi(x_i) = y_i\}$$

**Let's generalize it to Reward-Sensitive Classification:**

$$\{x_i, \mathbf{r}_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, \mathbf{r}_i \in [0, 1]^{|\mathcal{A}|}$$

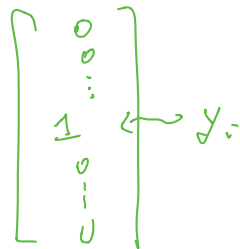
$\uparrow$   
feature

# Ingredient 2: Reward-sensitive Classification Oracle

Recall classic classification:

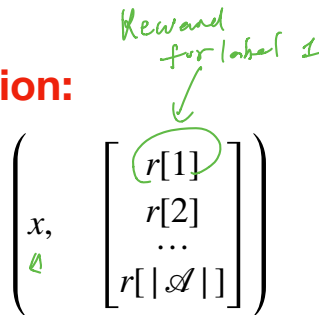
$$\{x_i, y_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, y_i \in \{1, \dots, |\mathcal{A}|\}$$

$$\arg \max_{\pi \in \Pi} \sum_{i=1}^N \mathbf{1}\{\pi(x_i) = y_i\}$$



Let's generalize it to **Reward-Sensitive Classification**:

$$\{x_i, \mathbf{r}_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, \mathbf{r}_i \in [0, 1]^{|\mathcal{A}|}$$





# Ingredient 2: Reward-sensitive Classification Oracle

Recall classic classification:

$$\{x_i, y_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, y_i \in \{1, \dots, |\mathcal{A}|\}$$

$$\arg \max_{\pi \in \Pi} \sum_{i=1}^N \mathbf{1}\{\pi(x_i) = y_i\}$$

Y[C:]

Let's generalize it to **Reward-Sensitive Classification**:

$$\{x_i, \mathbf{r}_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, \mathbf{r}_i \in [0, 1]^{|\mathcal{A}|}$$

$$\arg \max_{\pi \in \Pi} \sum_{i=1}^N \mathbf{r}_i[\pi(x_i)]$$

$$\left( x, \begin{bmatrix} r[1] \\ r[2] \\ \dots \\ r[|\mathcal{A}|] \end{bmatrix} \right)$$

# Ingredient 2: Reward-sensitive Classification Oracle

Recall classic classification:

$$\{x_i, y_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, y_i \in \{1, \dots, |\mathcal{A}|\}$$

$$\arg \max_{\pi \in \Pi} \sum_{i=1}^N \mathbf{1}\{\pi(x_i) = y_i\}$$

**Let's generalize it to Reward-Sensitive Classification:**

$$\{x_i, \mathbf{r}_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, \mathbf{r}_i \in [0, 1]^{|\mathcal{A}|}$$

$$\arg \max_{\pi \in \Pi} \sum_{i=1}^N \mathbf{r}_i[\pi(x_i)]$$

$$\left( x, \begin{bmatrix} r[1] \\ r[2] \\ \dots \\ r[|\mathcal{A}|] \end{bmatrix} \right)$$

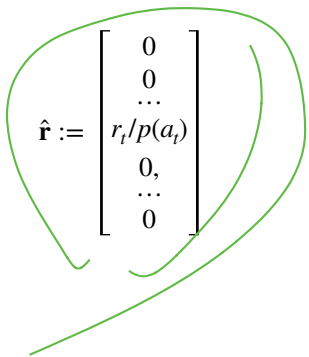
**We will do reduction to RSC**

# Summary so far:

1. Importance Weighting: we “magically” get unbiased estimate for all actions!

Assume  $a_t \sim p$  ( $p \in \Delta(\mathcal{A})$ ), For all  $a \in \mathcal{A}$ , define  $\hat{\mathbf{r}}[a] = \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)}$ , we have:

$$\mathbb{E}_{a_t \sim p} \hat{\mathbf{r}}[a] = r(x_t, a), \forall a \in \mathcal{A}$$


$$\hat{\mathbf{r}} := \begin{bmatrix} 0 \\ 0 \\ \dots \\ r_t/p(a_t) \\ 0, \\ \dots \\ 0 \end{bmatrix}$$

# Summary so far:

1. Importance Weighting: we “magically” get unbiased estimate for all actions!

Assume  $a_t \sim p$  ( $p \in \Delta(\mathcal{A})$ ), For all  $a \in \mathcal{A}$ , define  $\hat{\mathbf{r}}[a] = \frac{r(x_t, a)\mathbf{1}[a = a_t]}{p(a_t)}$ , we have:

$$\hat{\mathbf{r}} := \begin{bmatrix} 0 \\ 0 \\ \dots \\ r_t/p(a_t) \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

$$\mathbb{E}_{a_t \sim p} \hat{\mathbf{r}}[a] = r(x_t, a), \forall a \in \mathcal{A}$$

2. Reward-Sensitive Classification:

$$\{x_i, \mathbf{r}_i\}_{i=1}^N, \quad x_i \in \mathcal{X}, \mathbf{r}_i \in [0,1]^{|\mathcal{A}|}$$

$$\arg \max_{\pi \in \Pi} \sum_{i=1}^N \mathbf{r}_i[\pi(x_i)]$$

$$\left( x, \begin{bmatrix} r[1] \\ r[2] \\ \dots \\ r[|\mathcal{A}|] \end{bmatrix} \right)$$

# Put things together: Explore and Commit

For  $t = 0 \rightarrow N - 1$  : (# exploration phase)

# Put things together: Explore and Commit

For  $t = 0 \rightarrow N - 1$  : (# exploration phase)

1. Observe  $x_t \sim \mu$

# Put things together: Explore and Commit

For  $t = 0 \rightarrow N - 1$  : (# exploration phase)

1. Observe  $x_t \sim \mu$

uniform Dist over  $\mathcal{A}$

2. **Uniform-randomly** sample  $a_t \sim \text{Unif}(\mathcal{A})$ , receive reward  $r_t = r(x_t, a_t)$

# Put things together: Explore and Commit

For  $t = 0 \rightarrow N - 1$  : (# exploration phase)

1. Observe  $x_t \sim \mu$

2. **Uniform-randomly** sample  $a_t \sim \text{Unif}(\mathcal{A})$ , receive reward  $r_t = r(x_t, a_t)$

3. Use **IW**, form unbiased estimate  $\hat{\mathbf{r}}_t[a] = \begin{cases} 0 & a \neq a_t \\ \frac{r_t}{1/|\mathcal{A}|} & a = \underline{a_t} \end{cases}$

*All  $a \in \mathcal{A}$*



# Put things together: Explore and Commit

For  $t = 0 \rightarrow N - 1$  : (# exploration phase)

1. Observe  $x_t \sim \mu$

2. **Uniform-randomly** sample  $a_t \sim \text{Unif}(\mathcal{A})$ , receive reward  $r_t = r(x_t, a_t)$

3. Use **IW**, form unbiased estimate  $\hat{r}_t[a] = \begin{cases} 0 & a \neq a_t \\ \frac{r_t}{1/|\mathcal{A}|} & a = a_t \end{cases}$

$\{x_i, \hat{r}_i\}_{i=0}^{N-1} \leftarrow \text{Dataset for RSC}$   
 $\hat{r}_i \in \mathbb{R}^{|\mathcal{A}|}$

$$\hat{\mathbf{r}} := \begin{bmatrix} 0 \\ 0 \\ \dots \\ r_t/p(a_t) \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

$p(a_t) = \frac{1}{|\mathcal{A}|}$

# Put things together: Explore and Commit

For  $t = 0 \rightarrow N - 1$  : (# exploration phase)

1. Observe  $x_t \sim \mu$

2. **Uniform-randomly** sample  $a_t \sim \text{Unif}(\mathcal{A})$ , receive reward  $r_t = r(x_t, a_t)$

3. Use **IW**, form unbiased estimate  $\hat{\mathbf{r}}_t[a] = \begin{cases} 0 & a \neq a_t \\ \frac{r_t}{1/|\mathcal{A}|} & a = a_t \end{cases}$

$$\hat{\mathbf{r}} := \begin{bmatrix} 0 \\ 0 \\ \dots \\ r_t/p(a_t) \\ \dots \\ 0, \\ \dots \\ 0 \end{bmatrix}$$

**Call RSC oracle:**  $\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=0}^{N-1} \hat{\mathbf{r}}_i[\pi(x_i)]$

$x_t \sim \mu$

$\hat{\mathbf{r}}_t$  unbiased estimate of  $r(x_t, \cdot)$

for  $\pi$ :

$$\frac{1}{N} \sum_{i=1}^N \hat{Y}_i[\pi(x_i)]$$

is unbiased est of  $E_{x \sim \mu} [r(x, \pi(x))]$

$$\rightarrow \approx E_{x \sim \mu} [r(x, \pi(x))]$$

# Put things together: Explore and Commit

For  $t = 0 \rightarrow N - 1$  : (# exploration phase)

1. Observe  $x_t \sim \mu$

2. **Uniform-randomly** sample  $a_t \sim \text{Unif}(\mathcal{A})$ , receive reward  $r_t = r(x_t, a_t)$

3. Use **IW**, form unbiased estimate  $\hat{\mathbf{r}}_t[a] = \begin{cases} 0 & a \neq a_t \\ \frac{r_t}{1/|\mathcal{A}|} & a = a_t \end{cases}$

**Call RSC oracle:**  $\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=0}^{N-1} \hat{\mathbf{r}}_i[\pi(x_i)]$

$$E[\hat{F}_e] = r(x_t, \cdot)$$

$$\hat{\mathbf{r}} := \begin{bmatrix} 0 \\ 0 \\ \dots \\ r_t/p(a_t) \\ 0 \\ \dots \\ 0 \end{bmatrix}$$

For  $t = N \rightarrow T - 1$  : (# exploitation phase)

# Put things together: Explore and Commit

For  $t = 0 \rightarrow N - 1$  : (# exploration phase)

1. Observe  $x_t \sim \mu$

2. **Uniform-randomly** sample  $a_t \sim \text{Unif}(\mathcal{A})$ , receive reward  $r_t = r(x_t, a_t)$

3. Use **IW**, form unbiased estimate  $\hat{r}_t[a] = \begin{cases} 0 & a \neq a_t \\ \frac{r_t}{1/|\mathcal{A}|} & a = a_t \end{cases}$

**Call RSC oracle:**  $\hat{\pi} = \arg \max_{\pi \in \Pi} \sum_{i=0}^{N-1} \hat{r}_t[\pi(x_i)]$

$$\hat{\mathbf{r}} := \begin{bmatrix} 0 \\ 0 \\ \dots \\ r_t/p(a_t) \\ \dots \\ 0 \\ 0 \end{bmatrix}$$

For  $t = N \rightarrow T - 1$  : (# exploitation phase)

1. Observe  $x_t \sim \mu$ , and play  $a_t = \hat{\pi}(x_t)$

$$\frac{1}{N} \sum_{i=0}^{N-1} \hat{r}_t(\pi(x_i)) \Leftrightarrow E_{x \sim \mu} [r(x, \pi(x))]$$

## Outline for today:

 1. Introduction of the model

 2. Algorithm

3. Theory and some practical considerations

# Theory of the Explore and Commit Algorithm

For simplicity, assume  $\underline{\Pi}$  is discrete (but could be exponential large)

$\uparrow$  a class of classifiers

# Theory of the Explore and Commit Algorithm

For simplicity, assume  $\Pi$  is discrete (but could be exponential large)

[Theorem—informal] W/ high probability, properly setting the hyper-parameter  $N$ , Explore-and-Commit has the following regret:

$$\text{Regret}_T = T \mathbb{E}_{x \sim \mu}[r(x, \pi^*(x))] - \sum_{t=0}^{T-1} \mathbb{E}_{x \sim \mu}[r(x, \pi^t(x))] = O \left( \underbrace{T^{2/3} K^{1/3}}_{\text{Regret}_T} \cdot \underbrace{\ln(|\Pi|)^{1/3}}_{\substack{\uparrow \\ \text{Vc-Dim}(\Pi)}} \right)$$

$\frac{\text{Regret}_T}{T} \approx K^{\frac{1}{3}} T^{-\frac{1}{3}}$

# Practical Consideration

Instead of setting a hard threshold for explore and commit, we often interleave explore and exploitation



# Practical Consideration

Instead of setting a hard threshold for explore and commit, we often interleave explore and exploitation

$\epsilon \in (0, 1)$   
 $\epsilon$ -greedy:

Every iteration  $t$ :

With probability  $1 - \epsilon$ , we play  $a_t = \pi^t(x_t)$ ,

and w/ probability  $\epsilon$ , we play  $a_t \sim \text{Unif}(\mathcal{A})$

↖ Exploitation

↖ Exploration

# Practical Consideration

Instead of setting a hard threshold for explore and commit, we often interleave explore and exploitation

**$\epsilon$ -greedy:**

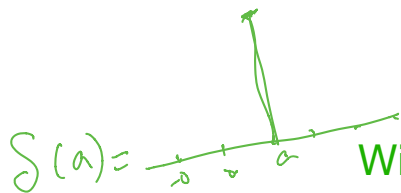
Every iteration  $t$ :

With probability  $1 - \epsilon$ , we play  $a_t = \pi^t(x_t)$ ,  
and w/ probability  $\epsilon$ , we play  $a_t \sim \text{Unif}(\mathcal{A})$

Q: What's the action distribution induced by  $\epsilon$ -greedy at iteration  $t$ ?

# Practical Consideration

Instead of setting a hard threshold for explore and commit, we often interleave explore and exploitation



**$\epsilon$ -greedy:**

Every iteration  $t$ :

With probability  $1 - \epsilon$ , we play  $a_t = \pi^t(x_t)$ ,  
and w/ probability  $\epsilon$ , we play  $a_t \sim \text{Unif}(\mathcal{A})$

Q: What's the action distribution induced by  $\epsilon$ -greedy at iteration  $t$ ?

$$a \sim p_t, \quad p_t = (1 - \epsilon)\delta(\pi^t(x_t)) + \epsilon \text{Unif}(\mathcal{A})$$

$\epsilon \rightarrow 0$  Exploitation

$\epsilon \rightarrow 1$  Uniform explore

# Put things together: $\epsilon$ -greedy

For  $t = 0 \rightarrow \infty$  (# interleave exploration & exploitation)

1. Observe  $x_t \sim \mu$

# Put things together: $\epsilon$ -greedy

For  $t = 0 \rightarrow \infty$  (# interleave exploration & exploitation)

1. Observe  $x_t \sim \mu$

2. Use  $\epsilon$ -greedy to form action distribution  $p_t = (1 - \epsilon)\delta(\pi^t(x_t)) + \epsilon\text{Unif}(\mathcal{A})$

# Put things together: $\epsilon$ -greedy

For  $t = 0 \rightarrow \infty$  (# interleave exploration & exploitation)

1. Observe  $x_t \sim \mu$
2. Use  $\epsilon$ -greedy to form action distribution  $p_t = (1 - \epsilon)\delta(\pi^t(x_t)) + \epsilon\text{Unif}(\mathcal{A})$
3. Play  $a_t \sim p_t$ , and observe  $r_t := r(x_t, a_t)$

# Put things together: $\epsilon$ -greedy

For  $t = 0 \rightarrow \infty$  (# interleave exploration & exploitation)

1. Observe  $x_t \sim \mu$

2. Use  $\epsilon$ -greedy to form action distribution  $p_t = (1 - \epsilon)\delta(\pi^t(x_t)) + \epsilon\text{Unif}(\mathcal{A})$

3. Play  $a_t \sim p_t$ , and observe  $r_t := r(x_t, a_t)$

4. Via IW, form unbiased estimate  $\hat{\mathbf{r}}_t$

A handwritten diagram in green ink showing a vertical vector structure. The vector is enclosed in large square brackets. From top to bottom, the elements are: 0, 0, a vertical ellipsis,  $\epsilon x_t / p_{\pi^t(x_t)}$ , 0, and  $\delta$ . A small arrow points to the top of the vector, and another small arrow points to the bottom element.

# Put things together: $\epsilon$ -greedy

For  $t = 0 \rightarrow \infty$  (# interleave exploration & exploitation)

1. Observe  $x_t \sim \mu$

2. Use  $\epsilon$ -greedy to form action distribution  $p_t = (1 - \epsilon)\delta(\pi^t(x_t)) + \epsilon\text{Unif}(\mathcal{A})$

3. Play  $a_t \sim p_t$ , and observe  $r_t := r(x_t, a_t)$

4. Via IW, form unbiased estimate  $\widehat{\mathbf{r}}_t$

5. Update via RSC oracle:  $\pi^{t+1} = \arg \max_{\pi \in \Pi} \sum_{i=1}^t \widehat{\mathbf{r}}_i[\pi(x_i)]$



# Put things together: $\epsilon$ -greedy


For  $t = 0 \rightarrow \infty$  (# interleave exploration & exploitation)

1. Observe  $x_t \sim \mu$

2. Use  $\epsilon$ -greedy to form action distribution  $p_t = (1 - \epsilon)\delta(\pi^t(x_t)) + \epsilon\text{Unif}(\mathcal{A})$

3. Play  $a_t \sim p_t$ , and observe  $r_t := r(x_t, a_t)$

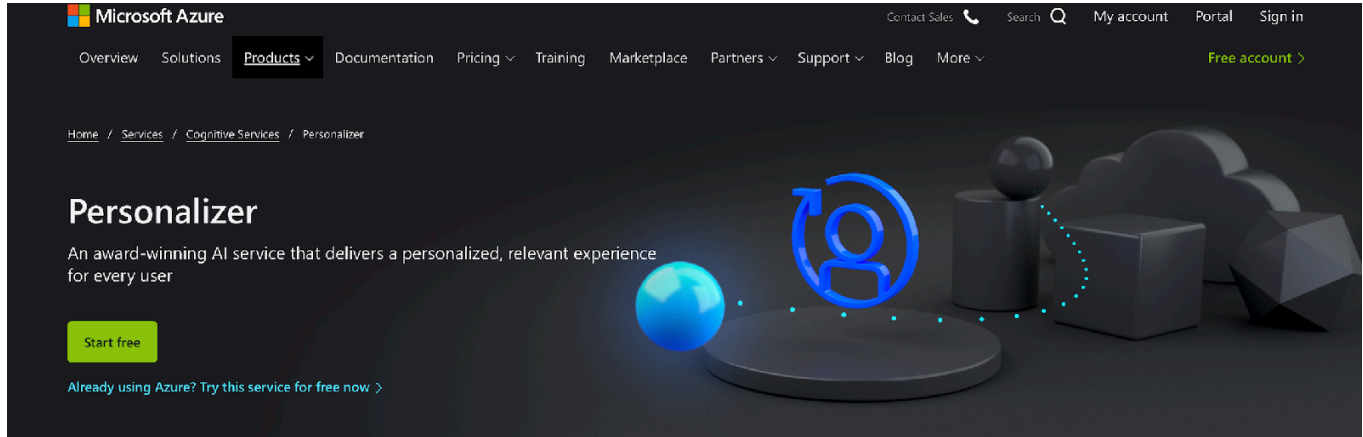
4. Via IW, form unbiased estimate  $\mathbf{r}_t$

5. Update via RSC oracle:  $\pi^{t+1} = \arg \max_{\pi \in \Pi} \sum_{i=1}^t \mathbf{r}_i[\pi(x_i)]$  

(Additionally 6. Gradually decay  $\epsilon \dots$ )

CB algorithm is being used in real world application at Microsoft:

<https://azure.microsoft.com/en-us/services/cognitive-services/personalizer/>



The image shows a screenshot of the Microsoft Azure website's product page for Personalizer. The page has a dark theme. At the top left is the Microsoft Azure logo. The top right contains navigation links: 'Contact Sales', 'Search', 'My account', 'Portal', and 'Sign in'. Below the logo is a horizontal menu with 'Overview', 'Solutions', 'Products' (highlighted with a dark background), 'Documentation', 'Pricing', 'Training', 'Marketplace', 'Partners', 'Support', and 'Blog'. A 'Free account >' link is on the far right. A breadcrumb trail reads 'Home / Services / Cognitive Services / Personalizer'. The main heading is 'Personalizer' in large white text. Below it is the tagline: 'An award-winning AI service that delivers a personalized, relevant experience for every user'. A green 'Start free' button is positioned to the left of a 3D graphic. The graphic features a blue sphere on a dark circular base, a blue '@' symbol with a circular arrow, and several dark grey geometric shapes (cylinder, cube, sphere) with a dotted blue line connecting them. At the bottom left, there is a link: 'Already using Azure? Try this service for free now >'. The background of the 3D graphic is dark with some light effects.

# Framework

