

# Online Stochastic Linear Optimization under One-bit Feedback

#### Lijun Zhang

#### LAMDA group, Nanjing University, China

The 6th Vision And Learning SEminar (VALSE 2016)

## Outline

### Introduction

- Definitions of Online Learning
- Bandit Online Learning
- 2 Online Learning under One-bit Feedback
  - Model and Algorithm
  - Theoretical Guarantees
  - Experimental Results





### What Happens in an Internet Minute?



http://cs.nju.edu.cn/zlj Linear Optimization under One-bit Feedback

## Outline

### Introduction

- Definitions of Online Learning
- Bandit Online Learning

### Online Learning under One-bit Feedback

- Model and Algorithm
- Theoretical Guarantees
- Experimental Results

### Conclusion and Future Work



・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・

### Online Algorithm vs Online Learning

#### Online Algorithm [Karp, 1992]

An online algorithm is one that receives a sequence of requests and performs an immediate action in response to each request.

- Computer Vision, Machine Learning, Data Mining
- Theoretical Computer Science, Computer Networks

#### Online Learning [Shalev-Shwartz, 2011]

Online learning is the process of answering a sequence of questions given (maybe partial) knowledge of <u>answers</u> to previous questions and possibly <u>additional information</u>.

Machine Learning, Game Theory, Information Theory



### Full-Information vs Bandit

#### Online Learning [Shalev-Shwartz, 2011]

Online learning is the process of <u>answering a sequence of</u> <u>questions</u> given (maybe partial) knowledge of <u>answers</u> to previous questions and possibly <u>additional information</u>.

• Full-Information Online Learning Multi-class Classification





### Full-Information vs Bandit

#### Online Learning [Shalev-Shwartz, 2011]

Online learning is the process of <u>answering a sequence of</u> <u>questions</u> given (maybe partial) knowledge of <u>answers</u> to previous questions and possibly <u>additional information</u>.

• Full-Information Online Learning Multi-class Classification



#### Bandit Online Learning



# **Formal Definitions**

#### **Online Learning**

1: for 
$$t = 1, 2, ..., T$$
 do

#### 4: end for



э

# Formal Definitions

### **Online Learning**

- 1: for t = 1, 2, ..., T do
- 2: Learner picks a decision  $\mathbf{x}_t \in \mathcal{D}$ Adversary chooses a function  $f_t(\cdot)$
- 4: end for



# **Formal Definitions**

### **Online Learning**

- 1: for t = 1, 2, ..., T do
- 2: Learner picks a decision  $\mathbf{x}_t \in \mathcal{D}$ Adversary chooses a function  $f_t(\cdot)$
- 3: Learner suffers loss  $f_t(\mathbf{x}_t)$
- 4: end for



# Formal Definitions

### **Online Learning**

- 1: for t = 1, 2, ..., T do
- 2: Learner picks a decision  $\mathbf{x}_t \in \mathcal{D}$ Adversary chooses a function  $f_t(\cdot)$
- 3: Learner suffers loss  $f_t(\mathbf{x}_t)$
- 4: end for

#### Cumulative Loss

Cumulative Loss = 
$$\sum_{t=1}^{T} \ell_t(\mathbf{x}_t)$$



# Formal Definitions

### **Online Learning**

- 1: for t = 1, 2, ..., T do
- 2: Learner picks a decision  $\mathbf{x}_t \in \mathcal{D}$ Adversary chooses a function  $f_t(\cdot)$
- 3: Learner suffers loss  $f_t(\mathbf{x}_t)$
- 4: end for

#### Regret

$$\text{Regret} = \sum_{t=1}^{T} \ell_t(\mathbf{x}_t) - \min_{\mathbf{x} \in \mathcal{D}} \sum_{t=1}^{T} f_t(\mathbf{x})$$



# **Formal Definitions**

### **Online Learning**

- 1: for t = 1, 2, ..., T do
- 2: Learner picks a decision  $\mathbf{x}_t \in \mathcal{D}$ Adversary chooses a function  $f_t(\cdot)$
- 3: Learner suffers loss  $f_t(\mathbf{x}_t)$
- 4: end for

#### Regret

$$\text{Regret} = \sum_{t=1}^{T} \ell_t(\mathbf{x}_t) - \min_{\mathbf{x} \in \mathcal{D}} \sum_{t=1}^{T} f_t(\mathbf{x})$$

- Full-Information Online Learning
  - Learner observes the function f<sub>t</sub>(·)
    [Zhang et al., 2012, Zhang et al., 2013]
- Bandit Online Learning
  - Learner only observes the value of  $f_t(\mathbf{x}_t)$

[Zhang et al., 2015, Zhang et al., 2016]

## Outline

### Introduction

- Definitions of Online Learning
- Bandit Online Learning

### 2 Online Learning under One-bit Feedback

- Model and Algorithm
- Theoretical Guarantees
- Experimental Results

### Conclusion and Future Work



## **Bandit Online Learning**

• Learner observes the value of  $f_t(\mathbf{x}_t)$  sequentially

Offline Counterpart: zero-order optimization



# **Bandit Online Learning**

- Learner observes the value of  $f_t(\mathbf{x}_t)$  sequentially
  - Offline Counterpart: zero-order optimization
- Learning Scenarios
  - Multi-Armed Bandits (MAB) [Robbins, 1952]
  - Multi-class Classification with Bandit Feedback
  - Online Convex Optimization with Bandit Feedback
    - Linear Bandits



不得下 イヨト イヨ

# Bandit Online Learning

- Learner observes the value of  $f_t(\mathbf{x}_t)$  sequentially
  - Offline Counterpart: zero-order optimization
- Learning Scenarios
  - Multi-Armed Bandits (MAB) [Robbins, 1952]
  - Multi-class Classification with Bandit Feedback
  - Online Convex Optimization with Bandit Feedback
    - Linear Bandits
- Generation Process of f<sub>t</sub>'s
  - Stochastic:  $f_1, \ldots, f_t$  are i.i.d.
  - Adversarial
    - Oblivious:  $f_t$  is independent of  $\mathbf{x}_1, \ldots, \mathbf{x}_{t-1}$  (like exam)
    - Nonoblivious:  $f_t$  depends on  $\mathbf{x}_1, \ldots, \mathbf{x}_{t-1}$  (like interview)



・ 同 ト ・ ヨ ト ・ ヨ ト

# Bandit Online Learning

- Learner observes the value of  $f_t(\mathbf{x}_t)$  sequentially
  - Offline Counterpart: zero-order optimization
- Learning Scenarios
  - Multi-Armed Bandits (MAB) [Robbins, 1952]
  - Multi-class Classification with Bandit Feedback
  - Online Convex Optimization with Bandit Feedback
    - Linear Bandits
- Generation Process of f<sub>t</sub>'s
  - Stochastic:  $f_1, \ldots, f_t$  are i.i.d.
  - Adversarial
    - Oblivious:  $f_t$  is independent of  $\mathbf{x}_1, \ldots, \mathbf{x}_{t-1}$  (like exam)
    - Nonoblivious:  $f_t$  depends on  $\mathbf{x}_1, \ldots, \mathbf{x}_{t-1}$  (like interview)



・ 同 ト ・ ヨ ト ・ ヨ ト

Multi-Armed Bandits (MAB)

• A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward





Multi-Armed Bandits (MAB)

- A gambler is facing K arms, and each time he pulls 1 arm and receives a reward Arm 1 Arm 2
  - Arm 3





Multi-Armed Bandits (MAB)

• A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward



( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( )



Multi-Armed Bandits (MAB)

• A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward



( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( )



Multi-Armed Bandits (MAB)

• A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward



A 3 1 A 3 1



Multi-Armed Bandits (MAB)

• A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward

Arm 1 $X_{1,1}$  $X_{1,2}$ Arm 210 $X_{2,2}$ Arm 3 $X_{3,1}$ 0



A 3 1 A 3 1



Multi-Armed Bandits (MAB)

• A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward





Multi-Armed Bandits (MAB)

• A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward



向下 イヨト イヨト



Multi-Armed Bandits (MAB)

- A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward



通り くほり くほう



Multi-Armed Bandits (MAB)

• A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward

Arm 1 $X_{1,1}$  $X_{1,2}$ 6 $X_{1,4}$ Arm 210 $X_{2,2}$  $X_{2,3}$ 0Arm 3 $X_{3,1}$ 0 $X_{3,3}$  $X_{3,4}$ 



通り くまり くまり



Multi-Armed Bandits (MAB)

• A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward



通り くまり くまり



Multi-Armed Bandits (MAB)

- A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward
  - Arm 1 $X_{1,1}$  $X_{1,2}$ 6 $X_{1,4}$  $X_{1,5}$ Arm 210 $X_{2,2}$  $X_{2,3}$ 0 $X_{2,5}$ Arm 3 $X_{3,1}$ 0 $X_{3,3}$  $X_{3,4}$  $X_{3,5}$
- Arm  $3 \mid X_{3,1} = 0$   $X_{3,3} = X_{3,4} = X$ • Exploration vs Exploitation



向下 イヨト イヨト



Multi-Armed Bandits (MAB)

- A gambler is facing *K* arms, and each time he pulls 1 arm and receives a reward
  - Arm 1 $X_{1,1}$  $X_{1,2}$ 6 $X_{1,4}$  $X_{1,5}$ Arm 210 $X_{2,2}$  $X_{2,3}$ 0 $X_{2,5}$ Arm 3 $X_{3,1}$ 0 $X_{3,3}$  $X_{3,4}$  $X_{3,5}$



3.5.4.3.5

- Exploration vs Exploitation
- Stochastic Setting
  - Rewards of the *i*-th arm are i.i.d. with unknown mean  $\mu_i$

$$\operatorname{Regret} = T \max_{i \in [K]} \mu_i - \sum_{t=1} \mu_{i_t}$$

• Upper Confidence Bound (UCB) [Auer et al., 2002]



### A Naive Approach based on Sample Mean





э

イロト イポト イヨト イヨト

### A Naive Approach based on Sample Mean





イロン イロン イヨン イヨン

### A Naive Approach based on Sample Mean





イロト イヨト イヨト

### A Naive Approach based on Sample Mean





イロト イポト イヨト イヨト



Definitions of Online Learning Bandit Online Learning

## Upper Confidence Bound (UCB)

### A Naive Approach based on Sample Mean


Introduction Learning under One-bit Feedback Conclusion

Definitions of Online Learning Bandit Online Learning

## Upper Confidence Bound (UCB)





3

<ロト < 四 > < 臣 > < 臣 > -

イロト イポト イヨト イヨト

## Upper Confidence Bound (UCB)

#### The Algorithm of UCB





э

イロン イロン イヨン イヨン

## Upper Confidence Bound (UCB)

### The Algorithm of UCB





イロト イヨト イヨト イヨト

## Upper Confidence Bound (UCB)

### The Algorithm of UCB





イロト イヨト イヨト イヨト

# Upper Confidence Bound (UCB)

### The Algorithm of UCB

With high probability  $\mu_i \leq \overline{\mu}_i = \hat{\mu}_i + \delta_i$ 





イロト イヨト イヨト イヨト

# Upper Confidence Bound (UCB)

### The Algorithm of UCB

With high probability  $\mu_i \leq \overline{\mu}_i = \hat{\mu}_i + \delta_i$ 





イロト イヨト イヨト イヨト

# Upper Confidence Bound (UCB)

### The Algorithm of UCB

With high probability  $\mu_i \leq \bar{\mu}_i = \hat{\mu}_i + \delta_i$ 





イロト イヨト イヨト イヨト

# Upper Confidence Bound (UCB)

### The Algorithm of UCB

With high probability  $\mu_i \leq \overline{\mu}_i = \hat{\mu}_i + \delta_i$ 





# Upper Confidence Bound (UCB)

### The Algorithm of UCB

With high probability  $\mu_i \leq \overline{\mu}_i = \hat{\mu}_i + \delta_i$ 



**Optimism in Face of Uncertainty**:  $i_t = \operatorname{argmax}_i \overline{\mu}_i$ 

• Construct  $\bar{\mu}_i$  by concentration inequalities (Chernoff-Hoeffding bound) [Auer et al., 2002]

$$\operatorname{Regret} = T \max_{i \in [K]} \mu_i - \sum_{t=1} \mu_{i_t} \le O(K \log T)$$

イロン イボン イヨン イヨ

# Outline

### Introduction

- Definitions of Online Learning
- Bandit Online Learning

Online Learning under One-bit Feedback

- Model and Algorithm
- Theoretical Guarantees
- Experimental Results
- Conclusion and Future Work



- Recommendation by Multi-Armed Bandits (MAB)
  - Each item is an arm
  - User feedback is the reward





- Recommendation by Multi-Armed Bandits (MAB)
  - Each item is an arm
  - User feedback is the reward





- Recommendation by Multi-Armed Bandits (MAB)
  - Each item is an arm
  - User feedback is the reward
  - The  $O(K \log T)$  regret bound is loose for large K
- Learning with Additional Information
  - Each arm is a feature vector  $\mathbf{x} \in \mathcal{D} \subseteq \mathbb{R}^d$







Model and Algorithm Theoretical Guarantees

- Recommendation by Multi-Armed Bandits (MAB)
  - Each item is an arm
  - User feedback is the reward
  - The  $O(K \log T)$  regret bound is loose for large K
- Learning with Additional Information
  - Each arm is a feature vector  $\mathbf{x} \in \mathcal{D} \subseteq \mathbb{R}^d$
  - For arm **x**, the expected reward  $\mu_{\mathbf{x}} = \mathbf{x}^{\top} \mathbf{w}_{*}$







Model and Algorithm Theoretical Guarantees Experiments

- Recommendation by Multi-Armed Bandits (MAB)
  - Each item is an arm
  - User feedback is the reward
  - The  $O(K \log T)$  regret bound is loose for large K
- Learning with Additional Information
  - Each arm is a feature vector  $\mathbf{x} \in \mathcal{D} \subseteq \mathbb{R}^d$
  - For arm **x**, the expected reward  $\mu_{\mathbf{x}} = \mathbf{x}^{\top} \mathbf{w}_{*}$

Regret = 
$$T \max_{\mathbf{x} \in \mathcal{D}} \mathbf{x}^\top \mathbf{w}_* - \sum_{t=1}^{T} \mathbf{x}_t^\top \mathbf{w}_*$$





Model and Algorithm Theoretical Guarantees Experiments

# Linear Optimization under One-bit Feedback

- Recommendation by Multi-Armed Bandits (MAB)
  - Each item is an arm
  - User feedback is the reward
  - The  $O(K \log T)$  regret bound is loose for large K
- Learning with Additional Information
  - Each arm is a feature vector  $\mathbf{x} \in \mathcal{D} \subseteq \mathbb{R}^d$
  - For arm **x**, the expected reward  $\mu_{\mathbf{x}} = \mathbf{x}^{\top} \mathbf{w}_{*}$

Regret = 
$$T \max_{\mathbf{x} \in \mathcal{D}} \mathbf{x}^\top \mathbf{w}_* - \sum_{t=1}^{T} \mathbf{x}_t^\top \mathbf{w}_*$$

Real-valued Feedback

[Dani et al., 2008]

$$\mathbf{y} = \mathbf{x}^\top \mathbf{w}_* + \epsilon \in \mathbb{R}$$







A (1) > A (2) > A (2)

Model and Algorithm Theoretical Guarantees

# Linear Optimization under One-bit Feedback

- Recommendation by Multi-Armed Bandits (MAB)
  - Each item is an arm
  - User feedback is the reward
  - The  $O(K \log T)$  regret bound is loose for large K
- Learning with Additional Information
  - Each arm is a feature vector  $\mathbf{x} \in \mathcal{D} \subseteq \mathbb{R}^d$
  - For arm **x**, the expected reward  $\mu_{\mathbf{x}} = \mathbf{x}^{\top} \mathbf{w}_{*}$

$$\operatorname{Regret} = T \max_{\mathbf{x} \in \mathcal{D}} \mathbf{x}^\top \mathbf{w}_* - \sum_{t=1}^{I} \mathbf{x}_t^\top \mathbf{w}_*$$

Real-valued Feedback [Dani et al., 2008]

 $\mathbf{y} = \mathbf{x}^\top \mathbf{w}_* + \epsilon \in \mathbb{R}$ 

**One-bit Feedback** [Zhang et al., 2016]

・部・・モー・ 不良に







Model and Algorithm Theoretical Guarantees

# Linear Optimization under One-bit Feedback

- Recommendation by Multi-Armed Bandits (MAB)
  - Each item is an arm
  - User feedback is the reward
  - The  $O(K \log T)$  regret bound is loose for large K
- Learning with Additional Information
  - Each arm is a feature vector  $\mathbf{x} \in \mathcal{D} \subseteq \mathbb{R}^d$
  - For arm **x**, the expected reward  $\mu_{\mathbf{x}} = \mathbf{x}^{\top} \mathbf{w}_{*}$

$$\operatorname{Regret} = T \max_{\mathbf{x} \in \mathcal{D}} \mathbf{x}^\top \mathbf{w}_* - \sum_{t=1}^{T} \mathbf{x}_t^\top \mathbf{w}_*$$

Real-valued Feedback [Dani et al., 2008]

**One-bit Feedback** [Zhang et al., 2016]

 $\Pr[y = \pm 1 | \mathbf{x}] = \frac{1}{1 + \exp(-y\mathbf{x}^{\top}\mathbf{w})}$ 

$$\mathbf{y} = \mathbf{x}^{\top} \mathbf{w}_* + \epsilon \in \mathbb{R}$$



Model and Algorithm Theoretical Guarantees Experiments

# Linear Optimization under One-bit Feedback

- Recommendation by Multi-Armed Bandits (MAB)
  - Each item is an arm
  - User feedback is the reward
  - The  $O(K \log T)$  regret bound is loose for large K
- Learning with Additional Information
  - Each arm is a feature vector  $\mathbf{x} \in \mathcal{D} \subseteq \mathbb{R}^d$
  - For arm **x**, the expected reward  $\mu_{\mathbf{x}} = \mathbf{x}^{\top} \mathbf{w}_{*}$

$$\operatorname{Regret} = T \max_{\mathbf{x} \in \mathcal{D}} \mathbf{x}^\top \mathbf{w}_* - \sum_{t=1}^{T} \mathbf{x}_t^\top \mathbf{w}_* \le O(d\sqrt{T})$$

Real-valued Feedback [Dani et al., 2008] One-bit Feedback [Zhang et al., 2016]

$$y = \mathbf{x}^{\top} \mathbf{w}_{*} + \epsilon \in \mathbb{R}$$
  $|\operatorname{Pr}[y = \pm 1 | \mathbf{x}] = \frac{1}{1 + \exp(-y\mathbf{x}^{\top}\mathbf{w})}$ 



イロト イポト イヨト イヨト

### A UCB-Type Algorithm [Zhang et al., 2016]

#### In the t-th Round

- Construct an upper bound  $\bar{\mu}_{\mathbf{x}} \ge \mu_{\mathbf{x}} = \mathbf{x}^{\top} \mathbf{w}_{*}$  for each arm  $\mathbf{x} \in \mathcal{D}$
- By Optimism in Face of Uncertainty,  $\mathbf{x}_t = \operatorname{argmax}_{\mathbf{x} \in D} \overline{\mu}_{\mathbf{x}}$



#### In the *t*-th Round



#### In the t-th Round



#### In the t-th Round



#### In the t-th Round



#### In the t-th Round



<ロ> (四) (四) (四) (日) (日)

### A UCB-Type Algorithm [Zhang et al., 2016]

#### In the t-th Round



<ロ> (四) (四) (四) (日) (日)

# A UCB-Type Algorithm [Zhang et al., 2016]

#### In the t-th Round

- Construct an upper bound  $\bar{\mu}_{\mathbf{x}} \ge \mu_{\mathbf{x}} = \mathbf{x}^{\top} \mathbf{w}_{*}$  for each arm  $\mathbf{x} \in \mathcal{D}$
- By Optimism in Face of Uncertainty,  $\mathbf{x}_t = \operatorname{argmax}_{\mathbf{x} \in \mathcal{D}} \overline{\mu}_{\mathbf{x}}$  $\mu_{\mathbf{x}} = \mathbf{x}^{\mathsf{T}} \mathbf{w}_{*} \leq \max_{\mathbf{w} \in \mathcal{C}_{t}} \mathbf{x}^{\mathsf{T}} \mathbf{w} = \bar{\mu}_{\mathbf{x}}$  $\mathbf{x}_t = \operatorname{argmax}_{\mathbf{x} \in \mathcal{D}} \overline{\mu}_{\mathbf{x}}$ Upper Bound Submit  $\mathbf{x}_t$ Observe  $y_t \in \{\pm 1\}$  $\mathbf{w}_* \in \mathcal{C}_t = \text{Ellipse}(\mathbf{w}_t, Z_t, \gamma_t)$ Save  $(\mathbf{x}_t, \mathbf{y}_t)$ **Confidence** Region  $\Pr[y_i = \pm 1] = \frac{1}{\exp(-y_i \mathbf{x}_i^{\mathsf{T}} \mathbf{w}_*)}$ Learning History  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_{t-1}, y_{t-1})$  $\mathbf{w}_t \in \mathbb{R}^d$

イロト イポト イヨト イヨ

## A UCB-Type Algorithm [Zhang et al., 2016]

#### In the t-th Round



## Implementation Issues [Zhang et al., 2016]

#### Online Newton Step

• Given 
$$\mathbf{w}_{t-1}$$
 and  $(\mathbf{x}_{t-1}, y_{t-1})$ 

• Define  $f_{t-1}(\mathbf{w}) = \log(1 + \exp(-y_{t-1}\mathbf{x}_{t-1}^{\top}\mathbf{w}))$ 

$$\begin{split} \mathbf{w}_t &= \operatorname*{argmin}_{\mathbf{w}} \frac{\|\mathbf{w} - \mathbf{w}_t\|_{Z_t}^2}{2} + (\mathbf{w} - \mathbf{w}_{t-1})^\top \nabla f_{t-1}(\mathbf{w}_{t-1}) \\ \text{where } Z_t &= Z_{t-1} + \frac{\beta}{2} \mathbf{x}_{t-1} \mathbf{x}_{t-1}^\top \end{split}$$

#### Arm Selection

$$\mathbf{X}_{t} = \underset{\mathbf{x} \in \mathcal{D}}{\operatorname{argmax}} \underbrace{\max_{\mathbf{w} \in \mathcal{C}_{t}} \mathbf{X}^{\top} \mathbf{w}}_{:=\overline{\mu}_{\mathbf{x}}} = \underset{\mathbf{x} \in \mathcal{D}}{\operatorname{argmax}} \max_{\|\mathbf{w} - \mathbf{w}_{t}\|_{Z_{t}} \leq \sqrt{\gamma_{t}}} \mathbf{X}^{\top} \mathbf{w}$$

- The above problem is NP-hard in general
- Tractable when  $\mathcal{D}$  is discrete or a ball



B 1 4 B 1

イロト イポト イヨト イヨ

## Outline

#### Introduction

- Definitions of Online Learning
- Bandit Online Learning

#### 2 Online Learning under One-bit Feedback

- Model and Algorithm
- Theoretical Guarantees
- Experimental Results

## Conclusion and Future Work



<ロト < 回 > < 回 > < 回 )

## Regret Bound [Zhang et al., 2016]

Theorem 1 (Confidence Region)

With a high probability, we have

$$(\mathbf{w}_* - \mathbf{w}_t)^{ op} \left( \lambda I + \frac{\beta}{2} \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^{ op} 
ight) (\mathbf{w}_* - \mathbf{w}_t) \le \gamma_t = O(d \log t)$$

for all t > 0.

- Optimality condition of online Newton step
- Bernstein's inequality for martingales
- Peeling technique for decoupling the dependence



### Regret Bound [Zhang et al., 2016]

Theorem 1 (Confidence Region)

With a high probability, we have

$$(\mathbf{w}_* - \mathbf{w}_t)^{ op} \left( \lambda I + \frac{\beta}{2} \sum_{i=1}^{t-1} \mathbf{x}_i \mathbf{x}_i^{ op} 
ight) (\mathbf{w}_* - \mathbf{w}_t) \leq \gamma_t = O(d \log t)$$

for all t > 0.

- Optimality condition of online Newton step
- Bernstein's inequality for martingales
- Peeling technique for decoupling the dependence

#### Theorem 2 (Regret)

With a high probability, we have

$$T \max_{\mathbf{x} \in \mathcal{D}} \mathbf{x}^{\top} \mathbf{w}_{*} - \sum_{t=1}^{T} \mathbf{x}_{t}^{\top} \mathbf{w}_{*} \leq 4 \sqrt{\frac{\gamma_{T} T}{\beta} \log \frac{\det(Z_{T+1})}{\det(Z_{1})}} = O\left(\frac{d\sqrt{T}}{D}\right)$$

for all T > 0.

イロト イポト イヨト イヨ

## Outline

#### Introduction

- Definitions of Online Learning
- Bandit Online Learning

### 2 Online Learning under One-bit Feedback

- Model and Algorithm
- Theoretical Guarantees
- Experimental Results

### 3 Conclusion and Future Work



# **Experimental Results**





## **Experimental Results**





# Experimental Results




 $\blacksquare \mathcal{D} \subseteq \mathbb{R}^{100} \text{ and } |\mathcal{D}| = 1000$ 



 $\blacksquare \mathcal{D} \subseteq \mathbb{R}^{100} \text{ and } |\mathcal{D}| = 1000$ 



$$\mathcal{D} = \{\mathbf{x} : \|\mathbf{x}\|_2 \le 1\} \subseteq \mathbb{R}^{100} \text{ and } |\mathcal{D}| = \infty$$





## **Conclusion and Future Work**

### Conclusion

- Online stochastic linear optimization under one-bit feedback
- An efficient algorithm with  $O(d\sqrt{T})$  regret bound

Bandit
Linear + One-bit
Full-Information

### Future Work

- Parameter selection in practice
- More observation models, such as the  $sign(\cdot)$
- Select multiple arms in each round
- w<sub>\*</sub> may change from round to round



### Resources



• Prediction, Learning, and Games [Cesa-Bianchi and Lugosi, 2006]







### Surveys

- Online Learning and Online Convex Optimization [Shalev-Shwartz, 2011]
- Regret Analysis of Stochastic and Nonstochastic Multi-armed Bandit Problems [Bubeck and Cesa-Bianchi, 2012]

### Reference I

# Thanks!



#### Auer, P., Cesa-Bianchi, N., and Fischer, P. (2002).

Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47(2-3):235–256.



#### Bubeck, S. and Cesa-Bianchi, N. (2012).

Regret analysis of stochastic and nonstochastic multi-armed bandit problems. Foundations and Trends in Machine Learning, 5(1):1–122.



Cesa-Bianchi, N. and Lugosi, G. (2006).

Prediction, Learning, and Games. Cambridge University Press.



Dani, V., Hayes, T. P., and Kakade, S. M. (2008).

Stochastic linear optimization under bandit feedback.

In Proceedings of the 21st Annual Conference on Learning, pages 355–366.



#### Karp, R. M. (1992).

On-line algorithms versus off-line algorithms: How much is it worth to know the future? In Proceedings of the IFIP 12th World Computer Congress on Algorithms, Software, Architecture – Information Processing, pages 416–429.



#### Robbins, H. (1952).

Some aspects of the sequential design of experiments. Bulletin of the American Mathematical Society, 58(5):527–535.



イロト イポト イヨト イヨト

### Reference II

#### 

#### Shalev-Shwartz, S. (2011).

Online learning and online convex optimization. Foundations and Trends in Machine Learning, 4(2):107–194.



Zhang, L., Jin, R., Chen, C., Bu, J., and He, X. (2012).

Efficient online learning for large-scale sparse kernel logistic regression. In Proceedings of the 26th AAAI Conference on Artificial Intelligence (AAAI), pages 1219–1225.



Zhang, L., Yang, T., Jin, R., Xiao, Y., and Zhou, Z.-H. (2016).

Online stochastic linear optimization under one-bit feedback.

In Proceedings of the 33rd International Conference on Machine Learning (ICML).



Zhang, L., Yang, T., Jin, R., and Zhou, Z.-H. (2015).

Online bandit learning for a special class of non-convex losses. In Proceedings of the 29th AAAI Conference on Artificial Intelligence (AAAI), pages 3158–3164.



Zhang, L., Yi, J., Jin, R., Lin, M., and He, X. (2013).

Online kernel learning with a near optimal sparsity bound. In Proceedings of the 30th International Conference on Machine Learning (ICML).



イロト イヨト イヨト イヨト