

# **Efficient Online Decision Tree Learning with Active Feature Acquisition**

**Arman Rahbar<sup>1</sup>, Ziyu Ye<sup>2</sup>, Yuxin Chen<sup>2</sup>, Morteza Haghiri Chehreghani<sup>1</sup>**

**<sup>1</sup>Chalmers University of Technology**

**<sup>2</sup>University of Chicago**

**IJCAI 2023**

# Agenda

- Introduction
- Proposed framework
- Results

# Introduction

Decision trees:

- Interpretable

Online variants of decision trees:

- Medical diagnosis
- Intrusion detection

Classical online DT learners:

- Require all features of incoming data points
- Not fully online

# Proposed framework

We take feature acquisition cost into account:

At each time step  $t$

1. Receive a data point with unknown features
2. Make prediction (i.e., classify) with low feature acquisition cost
3. Receive correct prediction for learning

Example: medical diagnosis:

At time  $t$ :

1. Patient  $\mathbf{x}^t$  comes in
  - $n$  medical tests:  $x_1^t, x_2^t, \dots, x_n^t$ , results unknown initially and can be measured at a cost
2. Predict accurate treatment with a low cost

# A New Problem Formulation

Data point:  $\mathbf{x} = (x_1, x_2, \dots, x_n)$

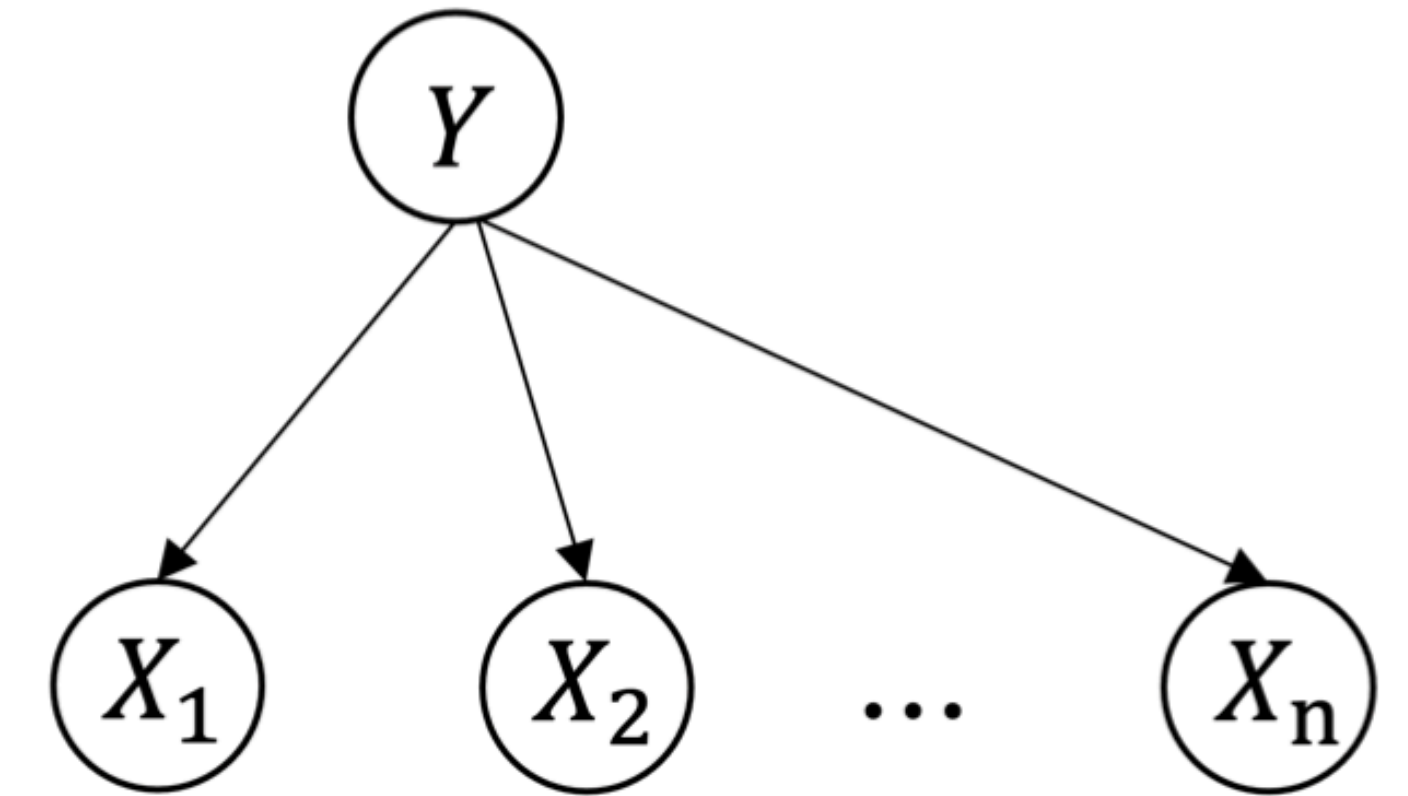
R.V. for feature values:  $X_i \in \mathcal{X} \triangleq \{0, 1\}$

Random variable for labels:  $Y_j \in \mathcal{Y} \triangleq \{y_1, y_2, \dots, y_m\}$

Naive Bayes assumption

- Joint distribution:  $\mathbb{P}[Y_j] \prod_{i=1}^n \mathbb{P}[X_i \mid Y_j]$

$$\theta_{ij} \triangleq \mathbb{P}[X_i = 1 \mid Y_j], \theta_{ij} \sim \text{Beta}(\alpha_{ij}, \beta_{ij})$$



# A New Problem Formulation

Random variable for *hypothesis* of data points:

full realization of features

$$H = [X_1, \dots, X_n], h \in \mathcal{H} \triangleq \{0, 1\}^n$$

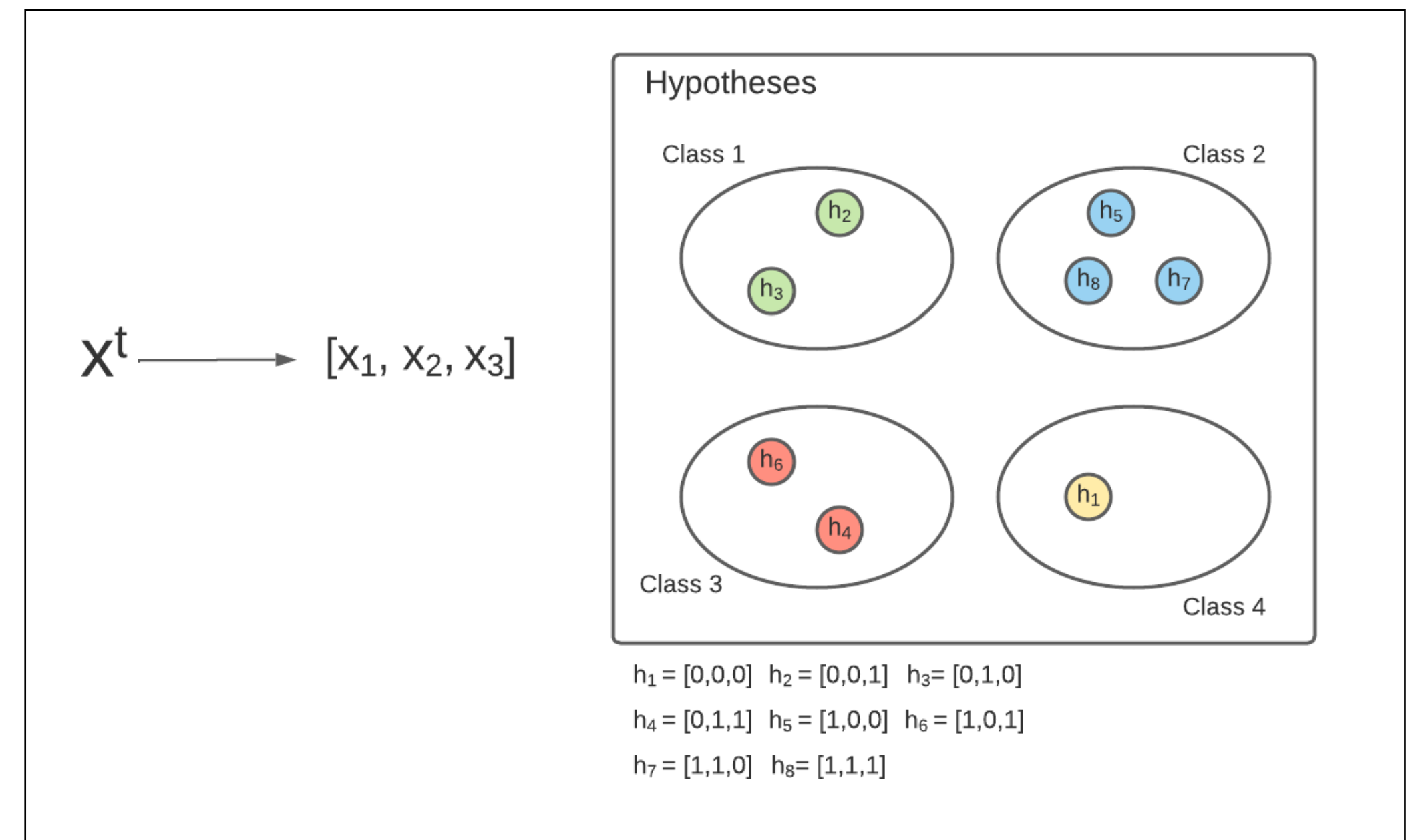
$\mathcal{H}$  partitioned into  $m$  disjoint

*decision regions / labels / classes*

Query  $q$  has cost  $c(q)$

prediction  $\hat{y}$  has loss  $l(\hat{y}, y)$

Goal: low  $c(q)$  and low  $l(\hat{y}, y)$



# Offline Prediction

$$\theta_{ij} \sim \text{Beta}(\alpha_{ij}, \beta_{ij})$$

$$\theta_{ij} \triangleq \mathbb{P}[X_i = 1 \mid Y_j]$$

Assume:

$\theta_{ij}$  for all  $ij$  and  $P(Y)$  are known

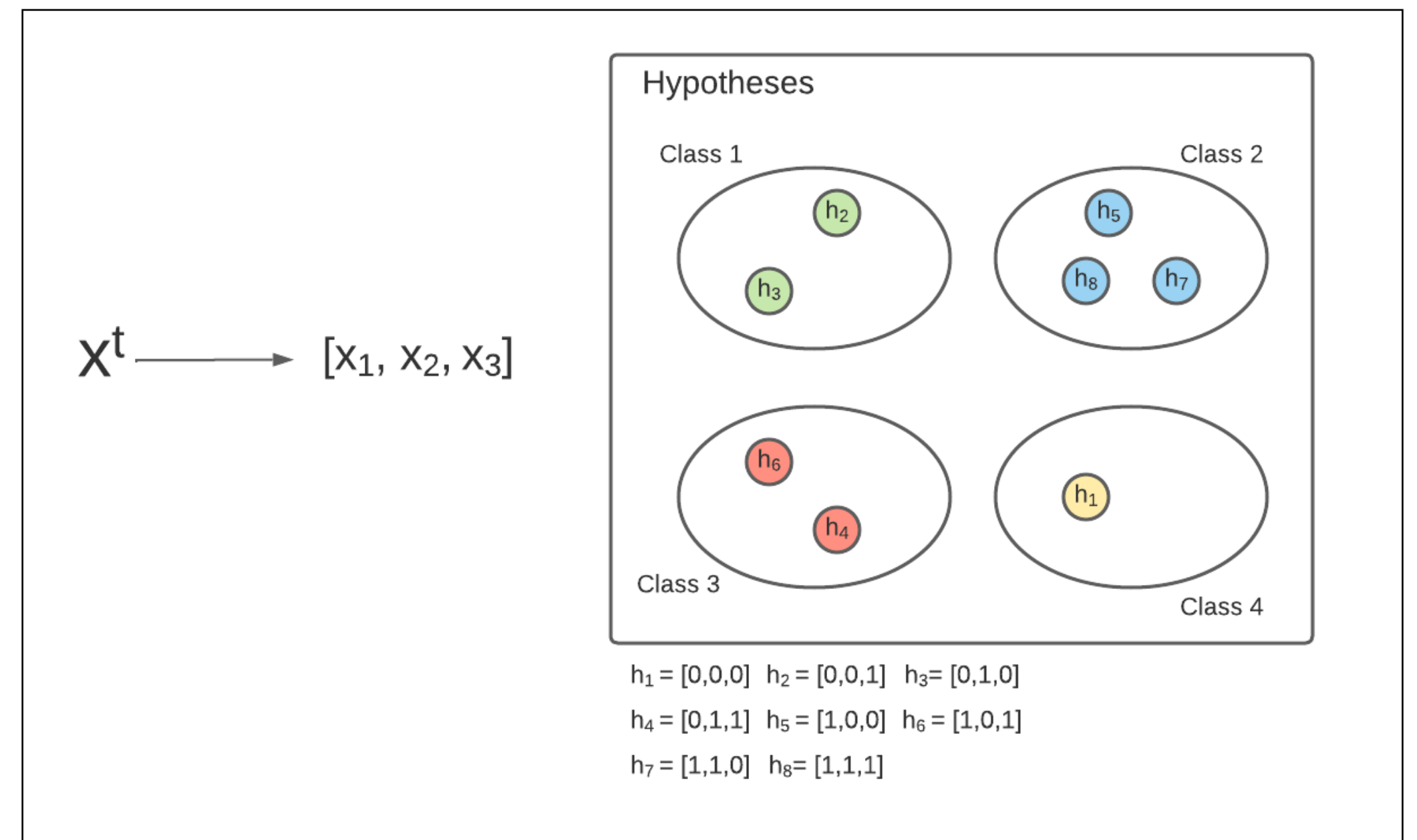
- We know the distribution of  $H$

Decision regions are known

Task: find decision region with low cost

Instance of:

Decision Region Determination (DRD)



# Decision Region Determination

Policy: mapping from current observations to features

Goal of DRD:

$$\pi^* = \arg \min_{\pi} \text{cost}(\pi)$$

s.t.  $\pi$  finds correct decision region

NP-hard

$EC^2$  [1] finds near-optimal policy and gives a sequence of features to query



# Online Prediction

We only have some prior knowledge (distributions):

- Use posterior sampling

# Online Prediction

# Online Prediction

1

# Online Prediction

**1 Sample the environment and receive data point**

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2**

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

Observe  $x_{\mathcal{F}}^t$



# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

Observe  $x_{\mathcal{F}}^t$

**3**

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

Observe  $x_{\mathcal{F}}^t$

**3 Make prediction and observe true label**

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

Observe  $x_{\mathcal{F}}^t$

**3 Make prediction and observe true label**

$$y_j^t$$

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

Observe  $x_{\mathcal{F}}^t$

**3 Make prediction and observe true label**

$y_j^t$

**4**

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

Observe  $x_{\mathcal{F}}^t$

**3 Make prediction and observe true label**

$y_j^t$

**4 Update the knowledge**

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

Observe  $x_{\mathcal{F}}^t$

**3 Make prediction and observe true label**

$y_j^t$

**4 Update the knowledge**

**for each**  $(i, x_i) \in x_{\mathcal{F}}^t$  **do**  
    **if**  $x_i = 1$  **then**  $\alpha_{ij}^t \leftarrow \alpha_{ij}^{t-1} + 1$   
    **else**  $\beta_{ij}^t \leftarrow \beta_{ij}^{t-1} + 1$

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

Observe  $x_{\mathcal{F}}^t$

**3 Make prediction and observe true label**

$y_j^t$

**4 Update the knowledge**

**for each**  $(i, x_i) \in x_{\mathcal{F}}^t$  **do**  
  **if**  $x_i = 1$  **then**  $\alpha_{ij}^t \leftarrow \alpha_{ij}^{t-1} + 1$   
  **else**  $\beta_{ij}^t \leftarrow \beta_{ij}^{t-1} + 1$

# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

Observe  $x_{\mathcal{F}}^t$

**3 Make prediction and observe true label**

$y_j^t$

**4 Update the knowledge**

**for each  $(i, x_i) \in x_{\mathcal{F}}^t$  do**  
    **if  $x_i = 1$  then  $\alpha_{ij}^t \leftarrow \alpha_{ij}^{t-1} + 1$**   
    **else  $\beta_{ij}^t \leftarrow \beta_{ij}^{t-1} + 1$**



# Online Prediction

**1 Sample the environment and receive data point**

$$\theta^t \sim \text{Beta}(\alpha^{t-1}, \beta^{t-1})$$

**2 Sample decision regions and use  $EC^2$  on the sampled environment**

Observe  $x_{\mathcal{F}}^t$

**3 Make prediction and observe true label**

$y_j^t$

**4 Update the knowledge**

Extensions:

- Real-valued features
- Concept drift

**for each**  $(i, x_i) \in x_{\mathcal{F}}^t$  **do**  
  **if**  $x_i = 1$  **then**  $\alpha_{ij}^t \leftarrow \alpha_{ij}^{t-1} + 1$   
  **else**  $\beta_{ij}^t \leftarrow \beta_{ij}^{t-1} + 1$

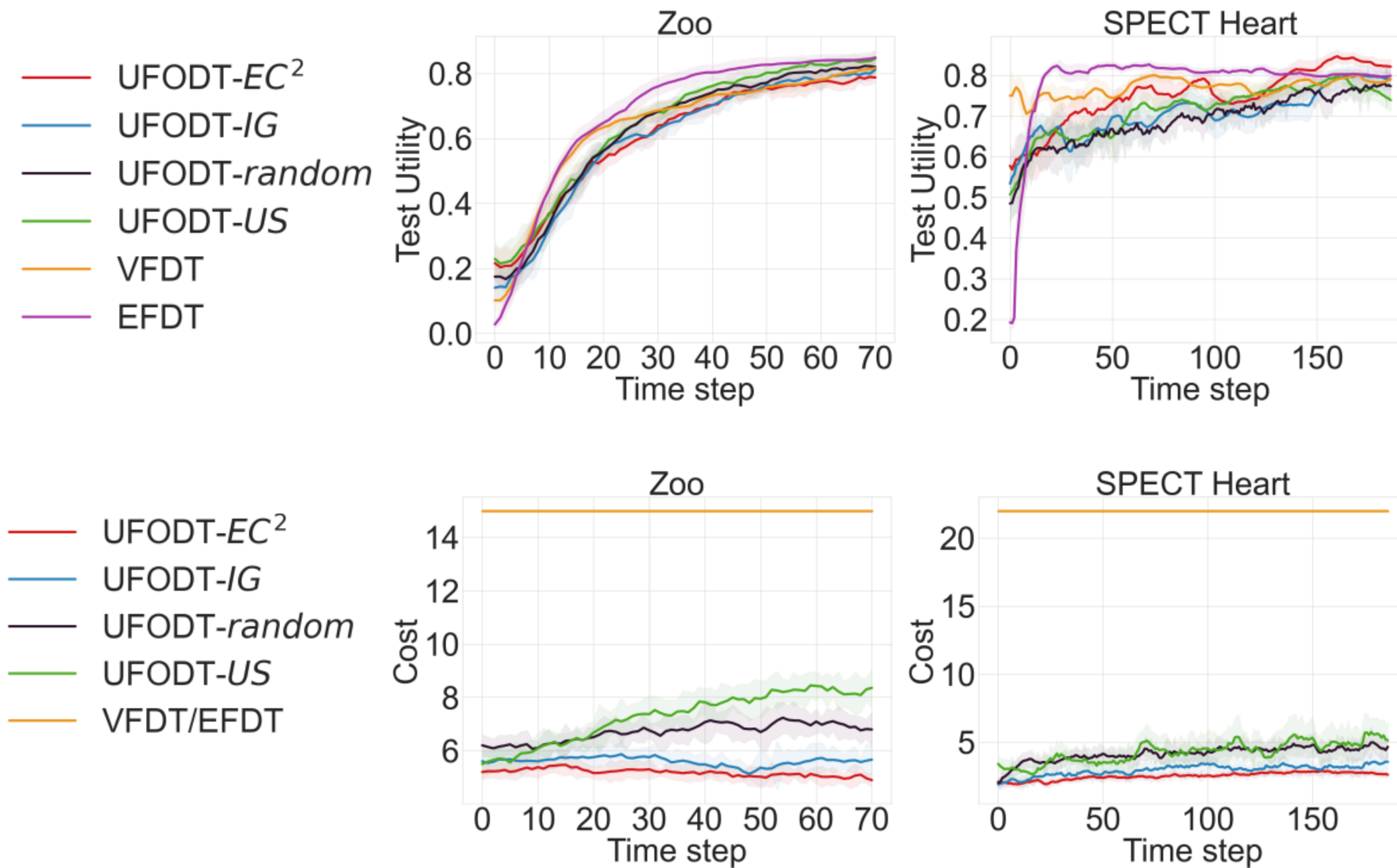
# Results

Theoretical analysis:

$$\Delta^t \triangleq \mathbb{U}(\pi_{\boldsymbol{\theta}^*}^*) - \mathbb{U}(\pi_{\boldsymbol{\theta}^t}^{\text{EC}^2})$$
$$\textit{Regret}(T) = \sum_{t=1}^T \Delta^t$$

Theorem:  $\mathbb{E}[\textit{Regret}(T)] = O(LS\sqrt{nLT \log(SnLT)})$

# Results



# References

[1] Daniel Golovin, Andreas Krause, and Debajyoti Ray. 2010. Near-optimal Bayesian active learning with noisy observations. In Proceedings of the 23rd International Conference on Neural Information Processing Systems - Volume 1 (NIPS'10). Curran Associates Inc., Red Hook, NY, USA, 766–774.