

Competitive Caching with Machine Learned Advice

Seminar on Online Algorithms

Gal Wiernik

Tel Aviv University

May 25, 2022

- 1 Introduction
- 2 Online Algorithms with ML Advice
- 3 The Predictive Marker Algorithm
- 4 Extensions
- 5 Experiments

Motivation - Machine Learning

Motivation - Machine Learning

- In recent years machine learning algorithms have been **wildly successful**.

Motivation - Machine Learning

- In recent years machine learning algorithms have been **wildly successful**.



Motivation - Machine Learning

- In recent years machine learning algorithms have been **wildly successful**.

- Yet, in practice:



Motivation - Machine Learning

- In recent years machine learning algorithms have been **wildly successful**.
- Yet, in practice:
 - Are very difficult to deploy



Motivation - Machine Learning

- In recent years machine learning algorithms have been **wildly successful**.
- Yet, in practice:
 - Are very difficult to deploy
 - Are prone to errors



Motivation - Machine Learning

- In recent years machine learning algorithms have been **wildly successful**.
- Yet, in practice:
 - Are very difficult to deploy
 - Are prone to errors



Motivation - Online Algorithms

- Online algorithms act without any knowledge of the future.

Motivation - Online Algorithms

- Online algorithms act without any knowledge of the future.
 - Are robust against any input

Motivation - Online Algorithms

- Online algorithms act without any knowledge of the future.
 - Are robust against any input
 - Have a provable guarantee on performance

Motivation - Online Algorithms

- Online algorithms act without any knowledge of the future.
 - Are robust against any input
 - Have a provable guarantee on performance
- Yet, overly cautious

Comparison

ML Algorithms

Online Algorithms

attempt to predict the unknown	act without any knowledge
suceptible to large errors	robust againt any input
exploit patterns	overly cautious

The main question

What if we could combine the **predictive power** of ML with the **robustness** of online algorithms?

First example

Example: Binary Search

Textbook problem - Sorted array A of size n , and query q .
What is the query cost?

First example

Example: Binary Search

Textbook problem - Sorted array A of size n , and query q .
What is the query cost?

ML Approach

First example

Example: Binary Search

Textbook problem - Sorted array A of size n , and query q .
What is the query cost?

ML Approach

- train a classifier $h(q)$ to predict $t(q)$.

First example

Example: Binary Search

Textbook problem - Sorted array A of size n , and query q .
What is the query cost?

ML Approach

- train a classifier $h(q)$ to predict $t(q)$.
- How can we use such a classifier?



$t(q)$



$h(q)$

$t(q)$

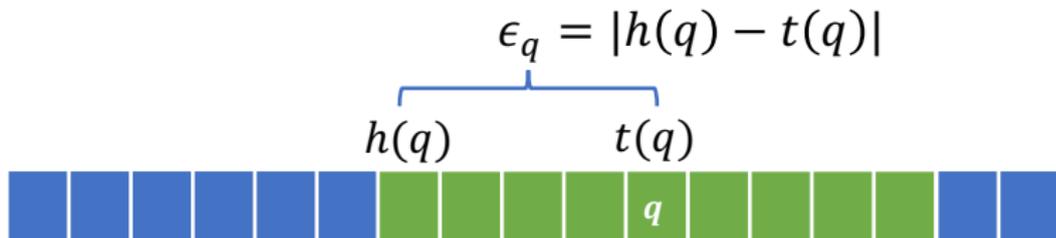




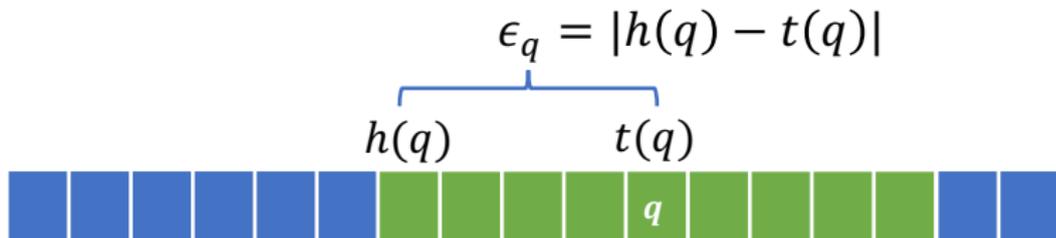
$$\epsilon_q = |h(q) - t(q)|$$

$h(q)$ $t(q)$





The expected cost is $2 \cdot \log(\epsilon_q)$.



The expected cost is $2 \cdot \log(\epsilon_q)$.
Is this any good?

The main result

The caching problem

- Our focus will be the **caching problem**.

The caching problem

- Our focus will be the **caching problem**.
- Best **deterministic** algorithm for online caching - $\Theta(k)$ competitive ratio

The caching problem

- Our focus will be the **caching problem**.
- Best **deterministic** algorithm for online caching - $\Theta(k)$ competitive ratio
- Best **randomized** algorithm - $\Theta(\log k)$

The caching problem

- Our focus will be the **caching problem**.
- Best **deterministic** algorithm for online caching - $\Theta(k)$ competitive ratio
- Best **randomized** algorithm - $\Theta(\log k)$
 - In reality, observed competitive ratio is much lower.

The caching problem

- Our focus will be the **caching problem**.
- Best **deterministic** algorithm for online caching - $\Theta(k)$ competitive ratio
- Best **randomized** algorithm - $\Theta(\log k)$
 - In reality, observed competitive ratio is much lower.

The main result

The **machine-learning assisted algorithm** reaches a competitive ratio of $2 + O(\min(\sqrt{\epsilon}, \log k))$.

- 1 Introduction
- 2 Online Algorithms with ML Advice
- 3 The Predictive Marker Algorithm
- 4 Extensions
- 5 Experiments

Preliminaries

- To achieve our results we have to define the playing ground for a **new genere of algorithms**: competitive algorithms with machine learning advice.

Preliminaries

- **ML scenarios** consist of:
 - Feature space - \mathcal{X} , and labels - \mathcal{Y}
 - hypothesis $h : \mathcal{X} \rightarrow \mathcal{Y}$
 - loss function: $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$

Preliminaries

- **ML scenarios** consist of:
 - Feature space - \mathcal{X} , and labels - \mathcal{Y}
 - hypothesis $h : \mathcal{X} \rightarrow \mathcal{Y}$
 - loss function: $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$
 - The loss can be: absolute ($\ell_1(y, \hat{y}) = |y - \hat{y}|$), squared ($\ell_2(y, \hat{y}) = (y - \hat{y})^2$), or generally: $\ell_c(y, \hat{y}) = \mathbf{1}_{y \neq \hat{y}}$

Preliminaries

- **ML scenarios** consist of:
 - Feature space - \mathcal{X} , and labels - \mathcal{Y}
 - hypothesis $h : \mathcal{X} \rightarrow \mathcal{Y}$
 - loss function: $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$
 - The loss can be: absolute ($\ell_1(y, \hat{y}) = |y - \hat{y}|$), squared ($\ell_2(y, \hat{y}) = (y - \hat{y})^2$), or generally: $\ell_c(y, \hat{y}) = \mathbf{1}_{y \neq \hat{y}}$
- **Online scenarios** consist of a algorithm \mathcal{A} and sequences σ .

Preliminaries

- **ML scenarios** consist of:
 - Feature space - \mathcal{X} , and labels - \mathcal{Y}
 - hypothesis $h : \mathcal{X} \rightarrow \mathcal{Y}$
 - loss function: $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$
 - The loss can be: absolute ($\ell_1(y, \hat{y}) = |y - \hat{y}|$), squared ($\ell_2(y, \hat{y}) = (y - \hat{y})^2$), or generally: $\ell_c(y, \hat{y}) = \mathbf{1}_{y \neq \hat{y}}$
- **Online scenarios** consist of a algorithm \mathcal{A} and sequences σ .
- \mathcal{A} has competitive ratio CR if for every σ :

$$\text{cost}_{\mathcal{A}}(\sigma) \leq \text{CR} \cdot \text{OPT}(\sigma)$$

The Online ML-Assisted Framework

- Now we can define the combined framework.

The Online ML-Assisted Framework

- Now we can define the combined framework.
- We have a **universe** \mathcal{Z} and **feature** space \mathcal{X}

The Online ML-Assisted Framework

- Now we can define the combined framework.
- We have a **universe** \mathcal{Z} and **feature** space \mathcal{X}
- The input is a sequence of items $\sigma = (\sigma_1, \sigma_2, \dots)$
Each item σ_i is associated with an element $z_i \in \mathcal{Z}$ and with features $x_i \in \mathcal{X}$

The Online ML-Assisted Framework

- Each item σ_i also has a label $y_i \in \mathcal{Y}$

The Online ML-Assisted Framework

- Each item σ_i also has a label $y_i \in \mathcal{Y}$
- A predictor $h : \mathcal{X} \rightarrow \mathcal{Y}$ returns $h(\sigma_i)$, attempts to find y_i

The Online ML-Assisted Framework

- Each item σ_i also has a label $y_i \in \mathcal{Y}$
- A predictor $h : \mathcal{X} \rightarrow \mathcal{Y}$ returns $h(\sigma_i)$, attempts to find y_i
 - The total loss of h on σ is:

$$\eta_\ell(h, \sigma) = \sum_i \ell(h(\sigma_i), y_i)$$

The Online ML-Assisted Framework

- Each item σ_i also has a label $y_i \in \mathcal{Y}$
- A predictor $h : \mathcal{X} \rightarrow \mathcal{Y}$ returns $h(\sigma_i)$, attempts to find y_i
 - The total loss of h on σ is:

$$\eta_\ell(h, \sigma) = \sum_i \ell(h(\sigma_i), y_i)$$

Question

How can we define h to have a general accuracy of ϵ ?

The Online ML-Assisted Framework

- Each item σ_i also has a label $y_i \in \mathcal{Y}$
- A predictor $h : \mathcal{X} \rightarrow \mathcal{Y}$ returns $h(\sigma_i)$, attempts to find y_i
 - The total loss of h on σ is:

$$\eta_\ell(h, \sigma) = \sum_i \ell(h(\sigma_i), y_i)$$

Question

How can we define h to have a general accuracy of ϵ ?

Definition

We say that h is ϵ -accurate if for every σ ,
 $\eta_\ell(h, \sigma) \leq \epsilon \cdot \text{OPT}(\sigma)$.

The Online ML-Assisted Framework

The Online ML-Assisted Framework

Define $CR_{\mathcal{A}}(\epsilon)$ to be the **competitive ratio** of algorithm \mathcal{A} that uses any predictor h that is ϵ -accurate.

The Online ML-Assisted Framework

Define $CR_{\mathcal{A}}(\epsilon)$ to be the **competitive ratio** of algorithm \mathcal{A} that uses any predictor h that is ϵ -accurate.



The Online ML-Assisted Framework

Define $CR_{\mathcal{A}}(\epsilon)$ to be the **competitive ratio** of algorithm \mathcal{A} that uses any predictor h that is ϵ -accurate.



- What would we like $CR_{\mathcal{A}}(0)$ to be?

The Online ML-Assisted Framework

Define $CR_{\mathcal{A}}(\epsilon)$ to be the **competitive ratio** of algorithm \mathcal{A} that uses any predictor h that is ϵ -accurate.



Consistency

\mathcal{A} is β -consistent if $CR_{\mathcal{A}}(0) = \beta$.

The Online ML-Assisted Framework

Define $CR_{\mathcal{A}}(\epsilon)$ to be the **competitive ratio** of algorithm \mathcal{A} that uses any predictor h that is ϵ -accurate.



Consistency

\mathcal{A} is β -consistent if $CR_{\mathcal{A}}(0) = \beta$.

Robustness

\mathcal{A} is α -robust for some function α if $CR_{\mathcal{A}}(\epsilon) = O(\alpha(\epsilon))$.

The Online ML-Assisted Framework

Define $CR_{\mathcal{A}}(\epsilon)$ to be the **competitive ratio** of algorithm \mathcal{A} that uses any predictor h that is ϵ -accurate.



Consistency

\mathcal{A} is β -consistent if $CR_{\mathcal{A}}(0) = \beta$.

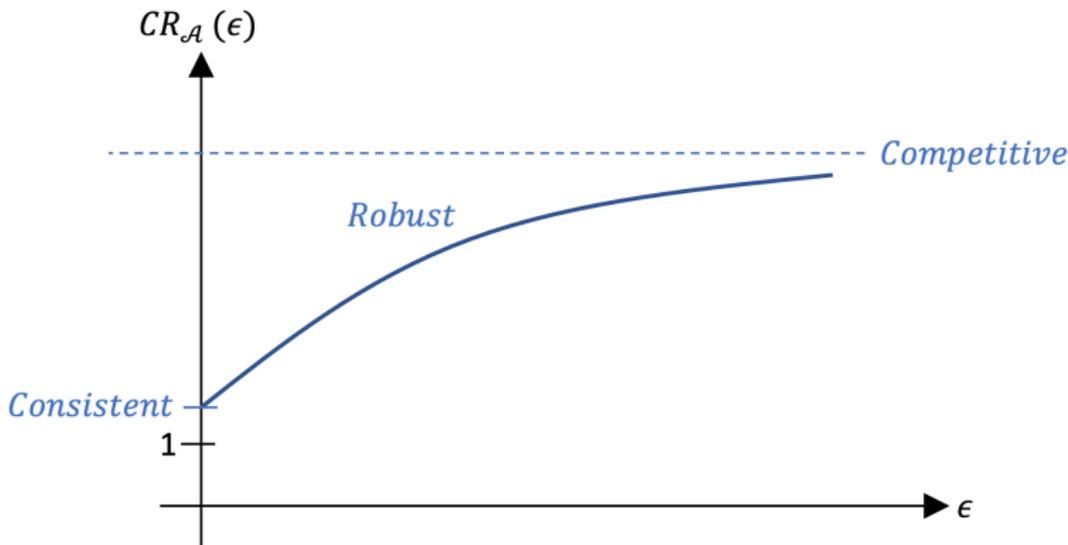
Robustness

\mathcal{A} is α -robust for some function α if $CR_{\mathcal{A}}(\epsilon) = O(\alpha(\epsilon))$.

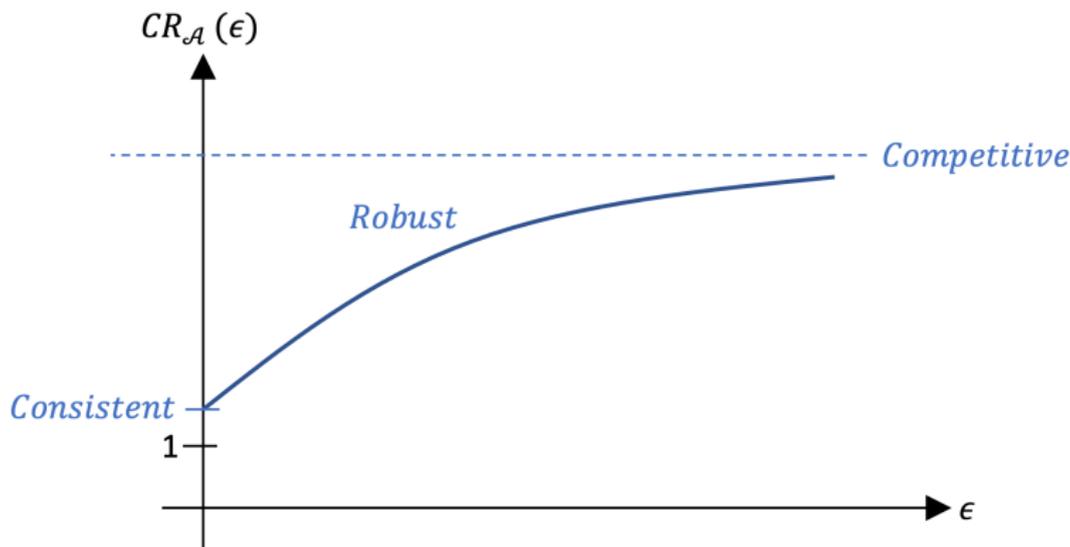
Competitiveness

\mathcal{A} is γ -competitive if $CR_{\mathcal{A}}(\epsilon) \leq \gamma \cdot \text{OPT}$ for all ϵ .

The Online ML-Assisted Framework



The Online ML-Assisted Framework



- The holy grail is an algorithm \mathcal{A} which simultaneously optimizes all three properties.

The Caching Scenario

The Caching Scenario

- Each σ_i is a request

The Caching Scenario

- Each σ_i is a request
- z_i is the requested page, and x_i is features of the request.

The Caching Scenario

- Each σ_i is a request
- z_i is the requested page, and x_i is features of the request.

Question

What should the labels \mathcal{Y} be?

Hint: The optimal caching algorithm is LFD.

The Caching Scenario

- Each σ_i is a request
- z_i is the requested page, and x_i is features of the request.

Question

What should the labels \mathcal{Y} be?

Hint: The optimal caching algorithm is LFD.

- $y(\sigma_i)$ will be the next appearance of z_i .

The Caching Scenario

- Each σ_i is a request
- z_i is the requested page, and x_i is features of the request.

Question

What should the labels \mathcal{Y} be?

Hint: The optimal caching algorithm is LFD.

- $y(\sigma_i)$ will be the next appearance of z_i .
- $h(\sigma_i)$ will try to predict y_i .

The Caching Scenario

- Each σ_i is a request
- z_i is the requested page, and x_i is features of the request.

Question

What should the labels \mathcal{Y} be?

Hint: The optimal caching algorithm is LFD.

- $y(\sigma_i)$ will be the next appearance of z_i .
- $h(\sigma_i)$ will try to predict y_i .
- $\mathcal{Y} = \mathbb{N}^+$

The Caching Scenario

- Example:

$y: a b c b a a c d \dots$

The diagram shows a sequence of characters $y: a b c b a a c d \dots$. Blue arrows indicate dependencies between elements: a top arrow from 'a' to 'b', a bottom arrow from 'b' to 'c', a bottom arrow from 'c' to 'b', a bottom arrow from 'b' to 'a', a bottom arrow from 'a' to 'a', and a bottom arrow from 'a' to 'c'.

The Caching Scenario

- Example:

$y: a b c b a a c d \dots$

The diagram shows a sequence of characters $y: a b c b a a c d \dots$. Blue arrows indicate dependencies between elements: a top arrow from 'a' to 'b', a bottom arrow from 'b' to 'c', a bottom arrow from 'c' to 'b', a bottom arrow from 'b' to 'a', a bottom arrow from 'a' to 'a', and a bottom arrow from 'a' to 'c'.

The Caching Scenario

- Example:

$y: a b c b a a c d \dots$

The diagram shows the sequence $y: a b c b a a c d \dots$. Blue arrows indicate dependencies: a top arrow from 'a' to 'b', a bottom arrow from 'c' to 'b', a bottom arrow from 'c' to 'a', and a bottom arrow from 'c' to 'a'.

$h: a b c b a a c d \dots$

The diagram shows the sequence $h: a b c b a a c d \dots$. Orange arrows indicate dependencies: a top arrow from 'a' to 'b', a bottom arrow from 'c' to 'b', a bottom arrow from 'c' to 'a', and a bottom arrow from 'c' to 'a'.

- 1 Introduction
- 2 Online Algorithms with ML Advice
- 3 The Predictive Marker Algorithm
- 4 Extensions
- 5 Experiments

Attempt #1 - Blindly Following the Predictor

- If the predictor is good, can't we just use it?

Attempt #1 - Blindly Following the Predictor

- If the predictor is good, can't we just use it?

Lemma

Define \mathcal{B} the algorithm that blindly follows an ϵ -accurate predictor. Then \mathcal{B} has a competitive ratio $\text{CR}_{\mathcal{B}}(\epsilon) = \Omega(\epsilon)$.

Attempt #1 - Blindly Following the Predictor

- If the predictor is good, can't we just use it?

Lemma

Define \mathcal{B} the algorithm that blindly follows an ϵ -accurate predictor. Then \mathcal{B} has a competitive ratio $\text{CR}_{\mathcal{B}}(\epsilon) = \Omega(\epsilon)$.

- Assume $k = 2$ and three elements, a, b, c .

Attempt #1 - Blindly Following the Predictor

- If the predictor is good, can't we just use it?

Lemma

Define \mathcal{B} the algorithm that blindly follows an ϵ -accurate predictor. Then \mathcal{B} has a competitive ratio $\text{CR}_{\mathcal{B}}(\epsilon) = \Omega(\epsilon)$.

- Assume $k = 2$ and three elements, a, b, c .
- Example follows.

a b c b c ... b c

Predictor h : *true, except $h(\sigma_1) = 2$*

a b c b c ... b c

Predictor h : *true, except $h(\sigma_1) = 2$*


 $a b c b c \dots b c$

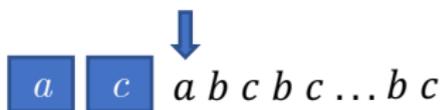
Predictor h : *true, except $h(\sigma_1) = 2$*

a b c b c ... b c

Predictor h : *true, except $h(\sigma_1) = 2$*

a c $a b c b c \dots b c$

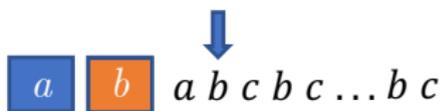
Predictor h : *true, except $h(\sigma_1) = 2$*

 a c $a b c b c \dots b c$

Predictor h : *true*, except $h(\sigma_1) = 2$

   $a b c b c \dots b c$

Predictor h : *true, except $h(\sigma_1) = 2$*

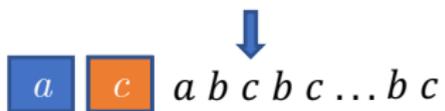
 a b $a b c b c \dots b c$

Predictor h : *true*, except $h(\sigma_1) = 2$

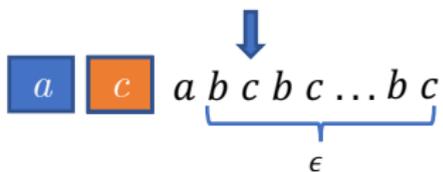
 $a b c b c \dots b c$

The diagram shows a sequence of characters. The first character 'a' is inside a blue square, and the second character 'b' is inside an orange square. To the right of these squares is the sequence 'a b c b c ... b c'. A blue arrow points downwards from the space between the 'a' and 'b' squares to the 'c' in the first 'a b c' triplet of the sequence.

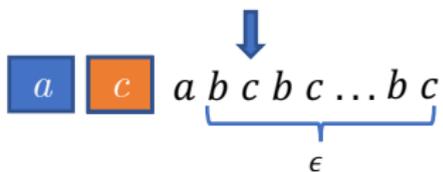
Predictor h : *true*, except $h(\sigma_1) = 2$

 a c $a b c b c \dots b c$

Predictor h : *true, except $h(\sigma_1) = 2$*

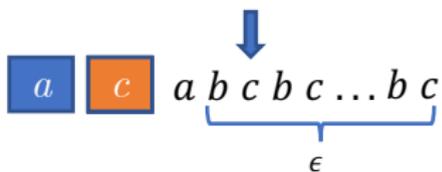


Predictor h : *true*, except $h(\sigma_1) = 2$



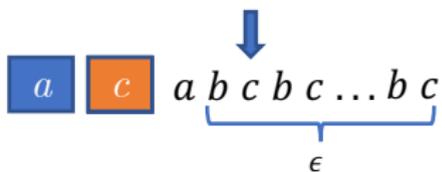
Offline Optimum:

Predictor h : true, except $h(\sigma_1) = 2$



Offline Optimum: $OPT = 1$

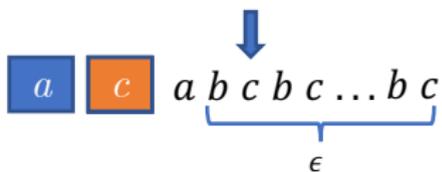
Predictor h : true, except $h(\sigma_1) = 2$



Offline Optimum: $OPT = 1$

Performance of \mathcal{B} :

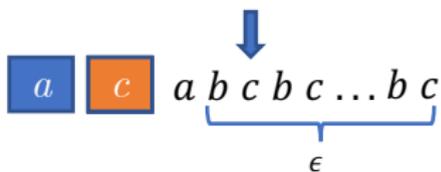
Predictor h : true, except $h(\sigma_1) = 2$



Offline Optimum: $OPT = 1$

Performance of \mathcal{B} : ϵ

Predictor h : true, except $h(\sigma_1) = 2$

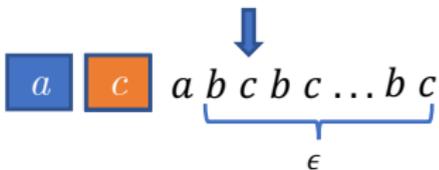


Offline Optimum: $OPT = 1$

Performance of \mathcal{B} : ϵ

Absolute Loss: $\eta(h, \sigma) = \epsilon$

Predictor h : true, except $h(\sigma_1) = 2$

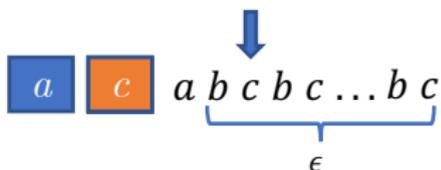


Offline Optimum: $OPT = 1$

Performance of \mathcal{B} : ϵ

Absolute Loss: $\eta(h, \sigma) = \epsilon \leq \epsilon \cdot OPT$

Predictor h : true, except $h(\sigma_1) = 2$



Offline Optimum: $OPT = 1$

Performance of \mathcal{B} : ϵ

Absolute Loss: $\eta(h, \sigma) = \epsilon \leq \epsilon \cdot OPT$

Competitive ratio: $CR_{\mathcal{B}}(\epsilon) = \frac{\epsilon}{1} = \Omega(\epsilon)$

Attempt #2 - Reacting to Predictor Mistakes

- We trusted the predictor too much. Can we do better?

Attempt #2 - Reacting to Predictor Mistakes

- We trusted the predictor too much. Can we do better?

Lemma

Define \mathcal{W} the algorithm that follows an ϵ -accurate predictor, but evicts wrong predictions. Then \mathcal{W} has a competitive ratio $\text{CR}_{\mathcal{W}}(\epsilon) = \Omega(\epsilon)$.

Attempt #2 - Reacting to Predictor Mistakes

- We trusted the predictor too much. Can we do better?

Lemma

Define \mathcal{W} the algorithm that follows an ϵ -accurate predictor, but evicts wrong predictions. Then \mathcal{W} has a competitive ratio $\text{CR}_{\mathcal{W}}(\epsilon) = \Omega(\epsilon)$.

- Assume $k = 3$ and four elements, a, b, c, d .

Attempt #2 - Reacting to Predictor Mistakes

- We trusted the predictor too much. Can we do better?

Lemma

Define \mathcal{W} the algorithm that follows an ϵ -accurate predictor, but evicts wrong predictions. Then \mathcal{W} has a competitive ratio $\text{CR}_{\mathcal{W}}(\epsilon) = \Omega(\epsilon)$.

- Assume $k = 3$ and four elements, a, b, c, d .
- Example follows.

d a b c a b c ... a b c

Predictor h : *predict d correctly*
but predict a, b, c two steps too early

$d a b c a b c \dots a b c$

Predictor h : *predict d correctly*
but predict a, b, c two steps too early


d a b c a b c ... a b c

Predictor h : *predict d correctly*
but predict a, b, c two steps too early

$d a b c a b c \dots a b c$

Predictor h : *predict d correctly*
but predict a, b, c two steps too early

a b c $d a b c a b c \dots a b c$

Predictor h : *predict d correctly*
but predict a, b, c two steps too early

 a b c d a b c a b c ... a b c

Predictor h : *predict d correctly*
but predict a, b, c two steps too early



Predictor h : *predict d correctly*
but predict a, b, c two steps too early

 a b d $d a b c a b c \dots a b c$

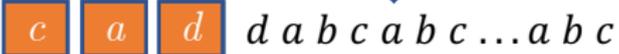
Predictor h : *predict d correctly*
but predict a, b, c two steps too early


 $d a b c a b c \dots a b c$

Predictor h : *predict d correctly*
but predict a, b, c two steps too early


 $d a b c a b c \dots a b c$

Predictor h : *predict d correctly*
but predict a, b, c two steps too early

Attempt #3 - Popular Heuristics

- In the examples we saw that there is some element that **should** have been evicted.

Attempt #3 - Popular Heuristics

- In the examples we saw that there is some element that **should** have been evicted.
- LRU and FIFO both provide strong heuristics for such cases.

Attempt #3 - Popular Heuristics

- In the examples we saw that there is some element that **should** have been evicted.
- LRU and FIFO both provide strong heuristics for such cases.
 - However, their **strict** (deterministic) policy leads to weak guarantees.

Reminder - Classic Marking Algorithm

- Recall the classic Marking Algorithm.

a b c c d c d a b b e a b

c c c

a b c c d c d a b b e a b

■ ■ ■ → a ■ ■ ■

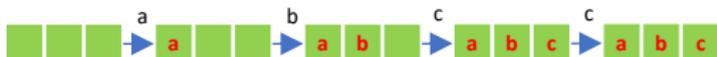
a b c c d c d a b b e a b



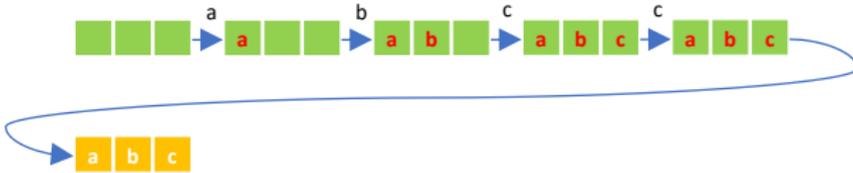
a b c c d c d a b b e a b



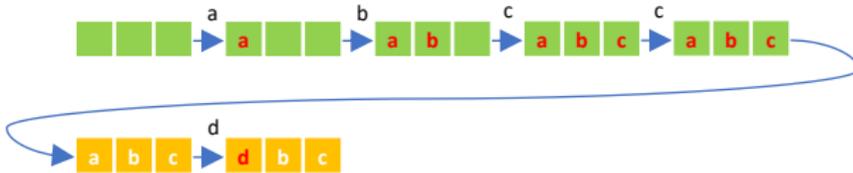
a b c c d c d a b b e a b



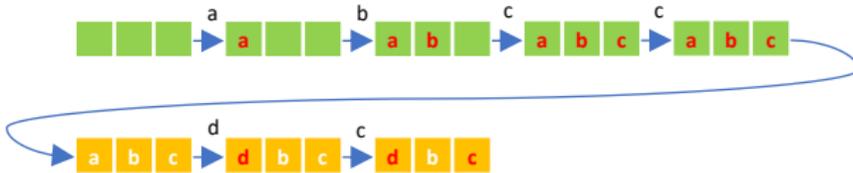
a b c c d c d a b b e a b



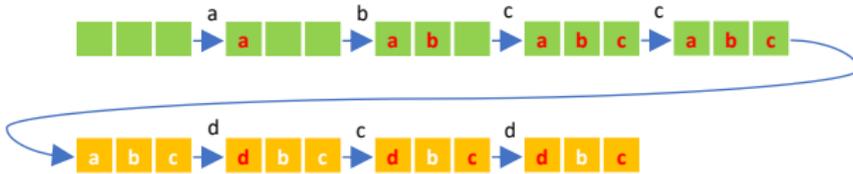
a b c c d c d a b b e a b



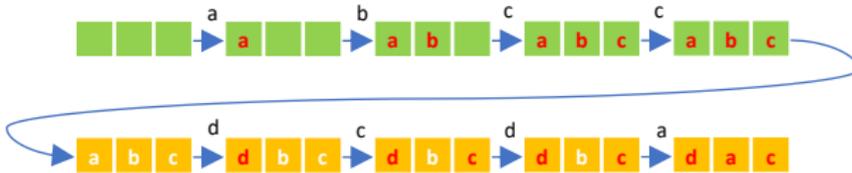
a b c c d c d a b b e a b



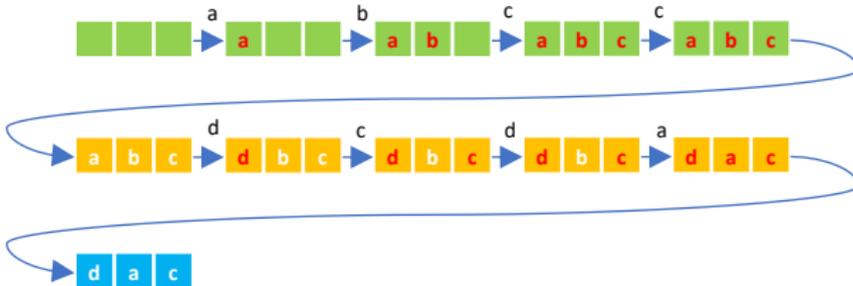
a b c c d c d a b b e a b



a b c c d c d a b b e a b



a b c c d c d a b b e a b



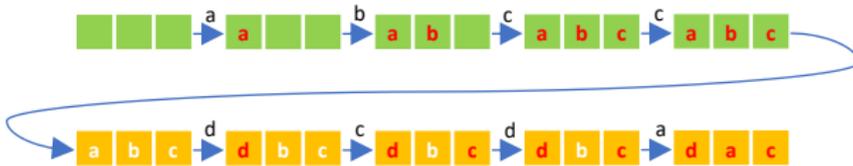
Reminder - Classic Marking Algorithm

- A **clean** element is an element that didn't arrive in phase $r - 1$ and arrives in phase r .

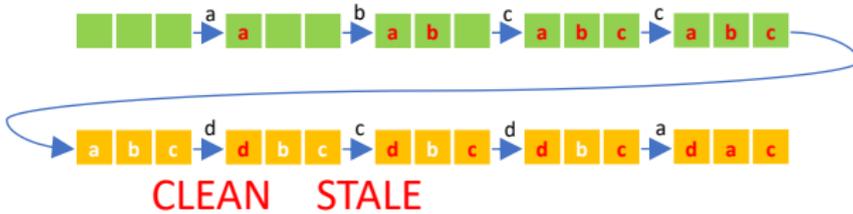
Reminder - Classic Marking Algorithm

- A **clean** element is an element that didn't arrive in phase $r - 1$ and arrives in phase r .
- A **stale** element is an element that arrived in phase $r - 1$ and **also** in phase r .

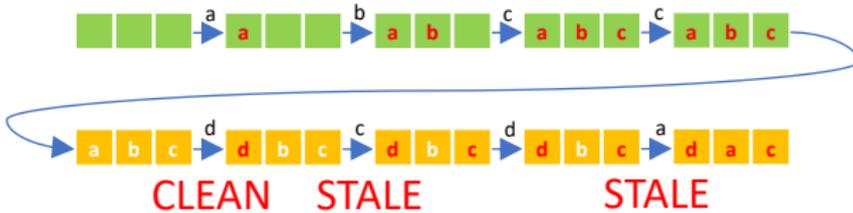
a b c c d c d a b b e a b



a b c c d c d a b b e a b



a b c c d c d a b b e a b



Reminder - Classic Marking Algorithm

- We saw that the marking algorithm has competitive ratio $O(\log k)$.

Reminder - Classic Marking Algorithm

- We saw that the marking algorithm has competitive ratio $O(\log k)$.
- This came from 2 claims:

Reminder - Classic Marking Algorithm

- We saw that the marking algorithm has competitive ratio $O(\log k)$.
- This came from 2 claims:

Claim 1

Let L be the number of clean elements in σ . Then $\text{OPT}(\sigma) \geq \frac{L}{2}$.

Reminder - Classic Marking Algorithm

- We saw that the marking algorithm has competitive ratio $O(\log k)$.
- This came from 2 claims:

Claim 1

Let L be the number of clean elements in σ . Then $\text{OPT}(\sigma) \geq \frac{L}{2}$.

Claim 2

Let L be the number of clean elements in σ .

Then $\mathbb{E}[\text{MARK}(\sigma)] \leq L \cdot H_k$.

Reminder - Classic Marking Algorithm

- We saw that the marking algorithm has competitive ratio $O(\log k)$.
- This came from 2 claims:

Claim 1

Let L be the number of clean elements in σ . Then $\text{OPT}(\sigma) \geq \frac{L}{2}$.

Claim 2

Let L be the number of clean elements in σ .

Then $\mathbb{E}[\text{MARK}(\sigma)] \leq L \cdot H_k$.

- Combined, we got:

$$\mathbb{E}[\text{MARK}(\sigma)] \leq L \cdot H_k \leq 2 \log k \cdot \text{OPT}(\sigma).$$

Predictive Marker

- If we'll use the marking algorithm, we'll gain $O(\log k)$ competitive ratio. How can we improve that?

Predictive Marker

- If we'll use the marking algorithm, we'll gain $O(\log k)$ competitive ratio. How can we improve that?
- The natural thing to do then is to **break ties in eviction** using the predictor.

Chains

Chains, Explained



Chains, Explained



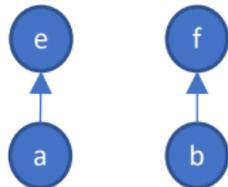
Chains, Explained



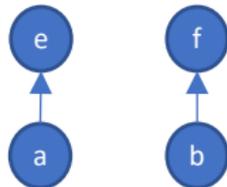
Chains, Explained



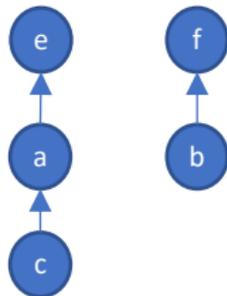
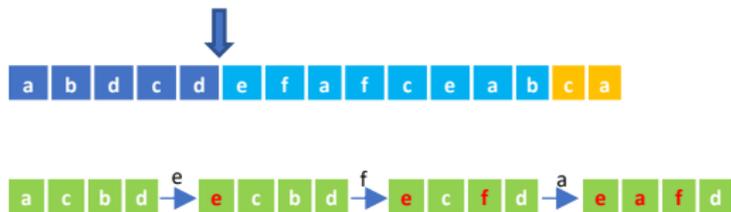
Chains, Explained



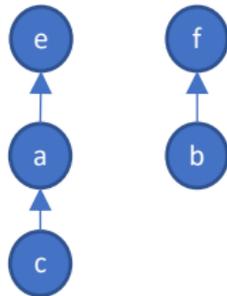
Chains, Explained



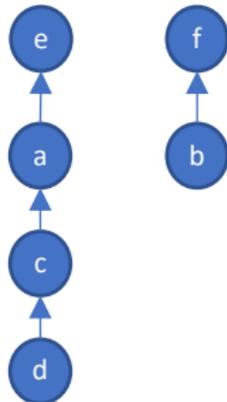
Chains, Explained



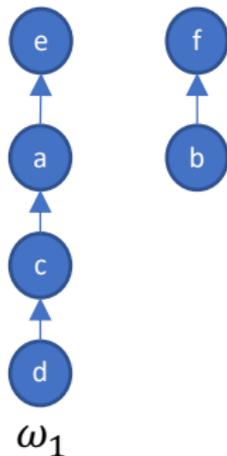
Chains, Explained



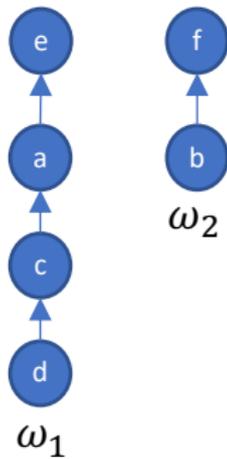
Chains, Explained



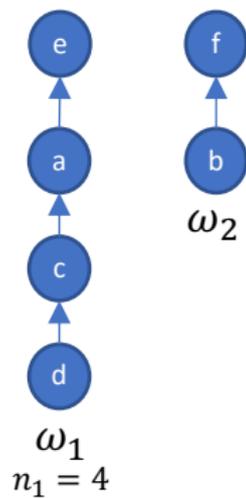
Chains, Explained



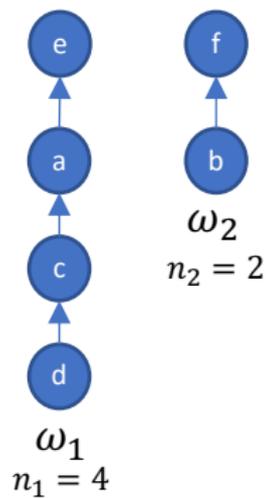
Chains, Explained



Chains, Explained



Chains, Explained



The Predictive Marker Algorithm

Now to the algorithm:

\mathcal{PM} – Predictive Marker

The Predictive Marker Algorithm

Now to the algorithm:

\mathcal{PM} – Predictive Marker

- At each phase, unmark all elements and save them as potentially **stale**.

The Predictive Marker Algorithm

Now to the algorithm:

\mathcal{PM} – Predictive Marker

- At each phase, unmark all elements and save them as potentially **stale**.
- In a cache miss for z_i :

The Predictive Marker Algorithm

Now to the algorithm:

\mathcal{PM} – Predictive Marker

- At each phase, unmark all elements and save them as potentially **stale**.
- In a cache miss for z_i :
 - If the element z_i is **clean**, create a new chain.

The Predictive Marker Algorithm

Now to the algorithm:

\mathcal{PM} – Predictive Marker

- At each phase, unmark all elements and save them as potentially **stale**.
- In a cache miss for z_i :
 - If the element z_i is **clean**, create a new chain.
 - If z_i is **stale**, find its chain $z_i = \omega_c$.

The Predictive Marker Algorithm

Now to the algorithm:

\mathcal{PM} – Predictive Marker

- At each phase, unmark all elements and save them as potentially **stale**.
- In a cache miss for z_i :
 - If the element z_i is **clean**, create a new chain.
 - If z_i is **stale**, find its chain $z_i = \omega_c$.
 - If the chain length is $n_c \leq H_k$:
evict unmarked element with highest predicted time e .

The Predictive Marker Algorithm

Now to the algorithm:

\mathcal{PM} – Predictive Marker

- At each phase, unmark all elements and save them as potentially **stale**.
- In a cache miss for z_i :
 - If the element z_i is **clean**, create a new chain.
 - If z_i is **stale**, find its chain $z_i = \omega_c$.
 - If the chain length is $n_c \leq H_k$:
 - evict unmarked element with highest predicted time e .
 - Else:
 - evict random unmarked element e .

The Predictive Marker Algorithm

Now to the algorithm:

\mathcal{PM} – Predictive Marker

- At each phase, unmark all elements and save them as potentially **stale**.
- In a cache miss for z_i :
 - If the element z_i is **clean**, create a new chain.
 - If z_i is **stale**, find its chain $z_i = \omega_c$.
 - If the chain length is $n_c \leq H_k$:
 - evict unmarked element with highest predicted time e .
 - Else:
 - evict random unmarked element e .
 - Increase the chain: $n_c \leftarrow n_c + 1$, $\omega_c \leftarrow e$.

Main Theorem

Theorem

Consider the caching scenario, with the prediction model \mathcal{H} and the loss function ℓ_1 .

The competitive ratio of the ϵ -assisted Predictive Marker Algorithm \mathcal{PM} is bounded by:

$$\text{CR}_{\mathcal{PM}}(\epsilon) \leq 2 \cdot \min\left(1 + \sqrt{5\epsilon}, 2H_k\right)$$

Proof - Chain Lengths

- Every cache miss is a link in a chain.

Proof - Chain Lengths

- Every cache miss is a link in a chain.
- Long chains means there are many misses. We would like to bound that length.

Proof - Chain Lengths

- Every cache miss is a link in a chain.
- Long chains means there are many misses. We would like to bound that length.

Lemma

If h has error $\leq \eta$ on chain ω_c , then the chain's length is bounded by $n_c \leq 1 + \sqrt{5\eta}$.

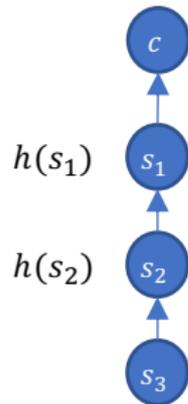
```
graph BT; s3((s3)) --> s2((s2)); s2 --> s1((s1)); s1 --> c((c))
```

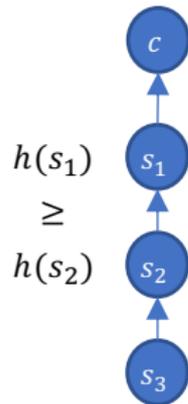
c

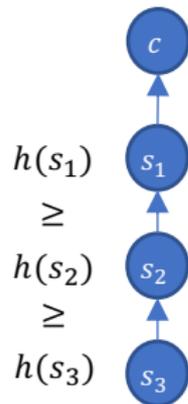
s_1

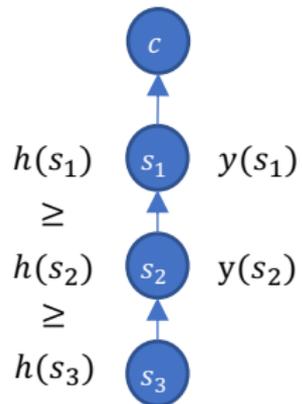
s_2

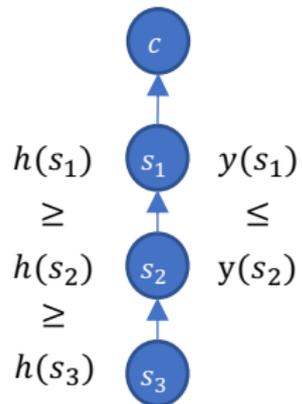
s_3

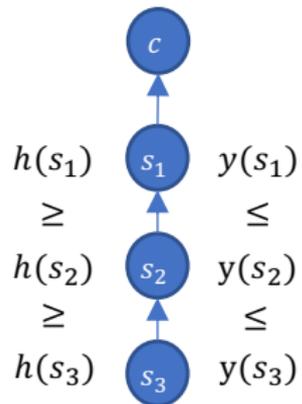


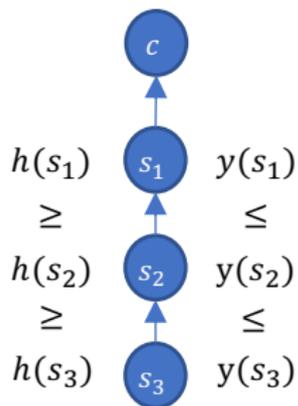








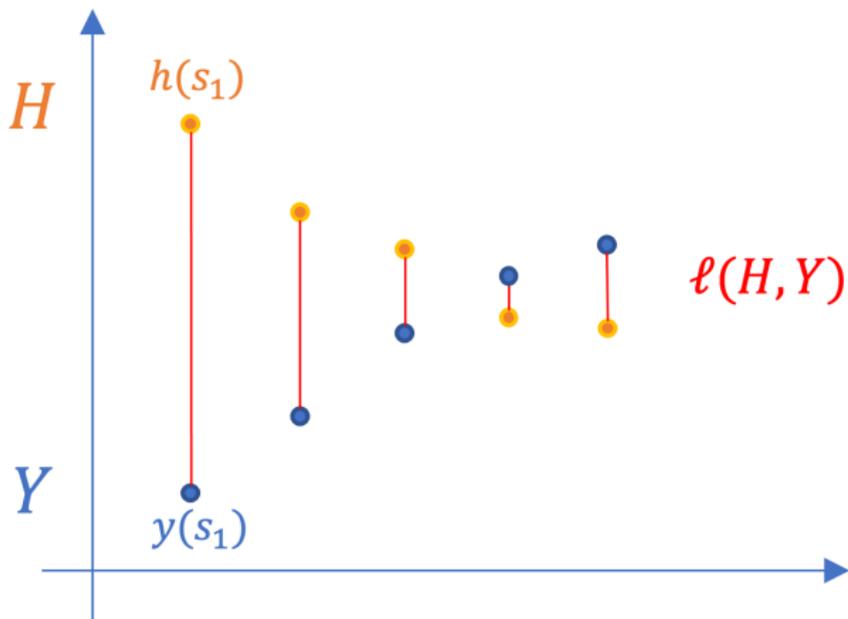




$$h(s_1) \geq h(s_2) \geq h(s_3) \geq \dots$$

$$y(s_1) \leq y(s_2) \leq y(s_3) \geq \dots$$

Proof - Chain Lengths

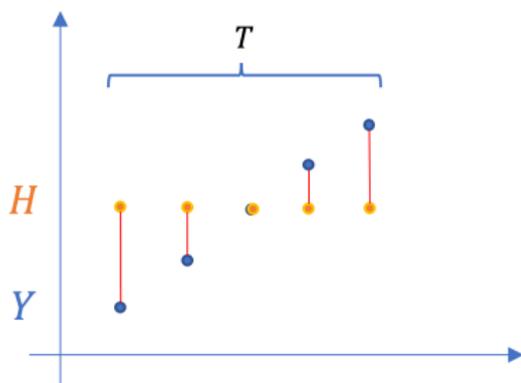


Proof - Chain Lengths

- If the error $\ell(H, Y)$ is bounded, how long can the chain get?

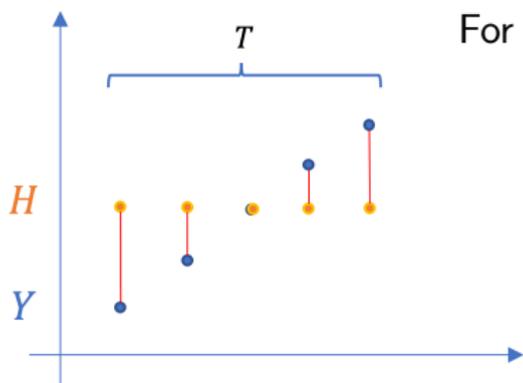
Proof - Chain Lengths

- If the error $\ell(H, Y)$ is bounded, how long can the chain get?



Proof - Chain Lengths

- If the error $\ell(H, Y)$ is bounded, how long can the chain get?

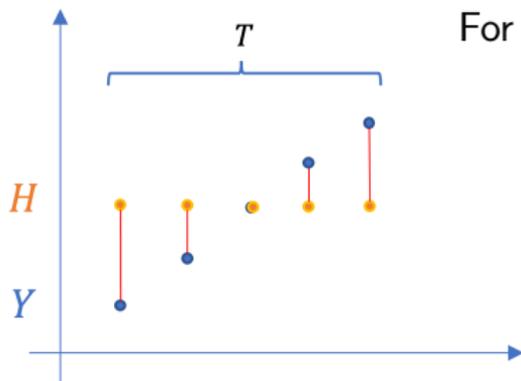


For such chain of length T we get the error:

$$\begin{aligned} \ell(H, Y) &= 2 \sum_{i=1}^{\frac{T-1}{2}} i = 2 \frac{\frac{T-1}{2} \cdot \frac{T+1}{2}}{2} \\ &= \frac{T^2 - 1}{4} \end{aligned}$$

Proof - Chain Lengths

- If the error $\ell(H, Y)$ is bounded, how long can the chain get?



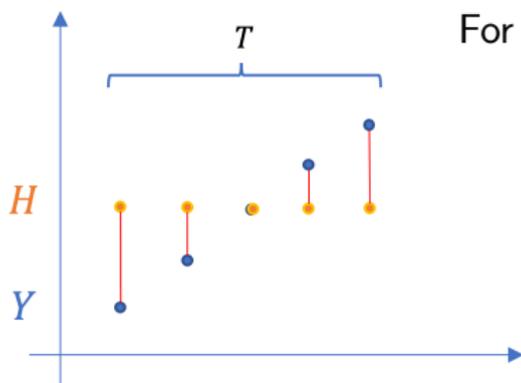
For such chain of length T we get the error:

$$\begin{aligned} \ell(H, Y) &= 2 \sum_{i=1}^{\frac{T-1}{2}} i = 2 \frac{\frac{T-1}{2} \cdot \frac{T+1}{2}}{2} \\ &= \frac{T^2 - 1}{4} \end{aligned}$$

- So every chain of length T has at least $\frac{T^2-1}{4}$ error.

Proof - Chain Lengths

- If the error $\ell(H, Y)$ is bounded, how long can the chain get?



For such chain of length T we get the error:

$$\begin{aligned} \ell(H, Y) &= 2 \sum_{i=1}^{T-1} i = 2 \frac{\frac{T-1}{2} \cdot \frac{T+1}{2}}{2} \\ &= \frac{T^2 - 1}{4} \end{aligned}$$

- So every chain of length T has at least $\frac{T^2-1}{4}$ error.

Corollary

If a chain has error $\leq \eta_c$, its length is at most $\sqrt{4\eta_c + 1} \leq \sqrt{5\eta_c}$.

Proof - Continued

- We want to figure out the total cost.

Proof - Continued

- We want to figure out the total cost.

Proof - Continued

- We want to figure out the total cost.
- For a single chain c with error η_c :

Proof - Continued

- We want to figure out the total cost.
- For a single chain c with error η_c :
 - If we never switched to random evictions: there are at most $1 + \sqrt{5\eta_c}$ cache misses.

Proof - Continued

- We want to figure out the total cost.
- For a single chain c with error η_c :
 - If we never switched to random evictions: there are at most $1 + \sqrt{5\eta_c}$ cache misses.
 - If we did: there are at most $2H_k$ misses.

Proof - Continued

- We want to figure out the total cost.
- For a single chain c with error η_c :
 - If we never switched to random evictions: there are at most $1 + \sqrt{5\eta_c}$ cache misses.
 - If we did: there are at most $2H_k$ misses.
- We can bound the evictions from the chain c of the r th phase in expectation by $\min(1 + \sqrt{5\eta_{r,c}}, 2H_k)$.

Proof - Continued

- We want to figure out the total cost.
- For a single chain c with error η_c :
 - If we never switched to random evictions: there are at most $1 + \sqrt{5\eta_c}$ cache misses.
 - If we did: there are at most $2H_k$ misses.
- We can bound the evictions from the chain c of the r th phase in expectation by $\min(1 + \sqrt{5\eta_{r,c}}, 2H_k)$.
- So the total cost is bounded by
$$\text{cost}_{\mathcal{P}, \mathcal{M}(\epsilon)} \leq \sum_{r,c} \min(1 + \sqrt{5\eta_{r,c}}, 2H_k) \leq ?$$

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).
The total error is $\eta = \epsilon \cdot \text{OPT}$.

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).

The total error is $\eta = \epsilon \cdot \text{OPT}$.

We want to know how big is $\sum_{r,c} \min(1 + \sqrt{5\eta_c}, 2H_k)$.

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).

The total error is $\eta = \epsilon \cdot \text{OPT}$.

We want to know how big is $\sum_{r,c} \min(1 + \sqrt{5\eta_c}, 2H_k)$.

- Since $\sqrt{\cdot}$ and $\min(\cdot)$ are both concave:

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).

The total error is $\eta = \epsilon \cdot \text{OPT}$.

We want to know how big is $\sum_{r,c} \min(1 + \sqrt{5\eta_c}, 2H_k)$.

- Since $\sqrt{\cdot}$ and $\min(\cdot)$ are both concave:
 - the way to maximize the expression is to divide the error η equally across all chains, with $\frac{\eta}{L}$ each.

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).

The total error is $\eta = \epsilon \cdot \text{OPT}$.

We want to know how big is $\sum_{r,c} \min(1 + \sqrt{5\eta_c}, 2H_k)$.

- Since $\sqrt{\cdot}$ and $\min(\cdot)$ are both concave:
 - the way to maximize the expression is to divide the error η equally across all chains, with $\frac{\eta}{L}$ each.
- The total length of all chains is then $L \cdot \min\left(1 + \sqrt{5\frac{\eta}{L}}, 2H_k\right)$.

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).

The total error is $\eta = \epsilon \cdot \text{OPT}$.

We want to know how big is $\sum_{r,c} \min(1 + \sqrt{5\eta_c}, 2H_k)$.

- Since $\sqrt{\cdot}$ and $\min(\cdot)$ are both concave:
 - the way to maximize the expression is to divide the error η equally across all chains, with $\frac{\eta}{L}$ each.
- The total length of all chains is then $L \cdot \min\left(1 + \sqrt{5\frac{\eta}{L}}, 2H_k\right)$.
- By Lemma 1, $\frac{L}{2} \leq \text{OPT}(\sigma)$. Trivially, $\text{OPT}(\sigma) \leq L$.

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).

The total error is $\eta = \epsilon \cdot \text{OPT}$.

We want to know how big is $\sum_{r,c} \min(1 + \sqrt{5\eta_c}, 2H_k)$.

- Since $\sqrt{\cdot}$ and $\min(\cdot)$ are both concave:
 - the way to maximize the expression is to divide the error η equally across all chains, with $\frac{\eta}{L}$ each.
- The total length of all chains is then $L \cdot \min\left(1 + \sqrt{5\frac{\eta}{L}}, 2H_k\right)$.
- By Lemma 1, $\frac{L}{2} \leq \text{OPT}(\sigma)$. Trivially, $\text{OPT}(\sigma) \leq L$.
- So: $\text{cost}_{\mathcal{P}\mathcal{M}(\epsilon)}(\sigma) \leq 2 \cdot \text{OPT}(\sigma) \cdot \min\left(1 + \sqrt{5\frac{\eta}{\text{OPT}(\sigma)}}, 2H_k\right)$.

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).

The total error is $\eta = \epsilon \cdot \text{OPT}$.

We want to know how big is $\sum_{r,c} \min(1 + \sqrt{5\eta_c}, 2H_k)$.

- Since $\sqrt{\cdot}$ and $\min(\cdot)$ are both concave:
 - the way to maximize the expression is to divide the error η equally across all chains, with $\frac{\eta}{L}$ each.
- The total length of all chains is then $L \cdot \min\left(1 + \sqrt{5\frac{\eta}{L}}, 2H_k\right)$.
- By Lemma 1, $\frac{L}{2} \leq \text{OPT}(\sigma)$. Trivially, $\text{OPT}(\sigma) \leq L$.
- So: $\text{cost}_{\mathcal{P}\mathcal{M}(\epsilon)}(\sigma) \leq 2 \cdot \text{OPT}(\sigma) \cdot \min\left(1 + \sqrt{5\frac{\eta}{\text{OPT}(\sigma)}}, 2H_k\right)$.

Proof - Continued

Proof.

Let L be the number of clean elements (= number of chains).

The total error is $\eta = \epsilon \cdot \text{OPT}$.

We want to know how big is $\sum_{r,c} \min(1 + \sqrt{5\eta_c}, 2H_k)$.

- Since $\sqrt{\cdot}$ and $\min(\cdot)$ are both concave:
 - the way to maximize the expression is to divide the error η equally across all chains, with $\frac{\eta}{L}$ each.
- The total length of all chains is then $L \cdot \min\left(1 + \sqrt{5\frac{\eta}{L}}, 2H_k\right)$.
- By Lemma 1, $\frac{L}{2} \leq \text{OPT}(\sigma)$. Trivially, $\text{OPT}(\sigma) \leq L$.
- So: $\text{cost}_{\mathcal{P}\mathcal{M}(\epsilon)}(\sigma) \leq 2 \cdot \text{OPT}(\sigma) \cdot \min\left(1 + \sqrt{5\frac{\eta}{\text{OPT}(\sigma)}}, 2H_k\right)$.

Which means $\text{CR}_{\mathcal{P}\mathcal{M}}(\epsilon) \leq 2 \cdot \min\left(1 + \sqrt{5\epsilon}, 2 \log k\right)$. □

Tightness of analysis

Theorem (without proof)

Tightness of analysis

Theorem (without proof)

Any **deterministic** ϵ -assisted marking algorithm \mathcal{A} , that only uses the predictor in tie-breaking among unmarked elements in a deterministic fashion, has a competitive ratio of

$$CR_{\mathcal{A}}(\epsilon) = \Omega(\min(\sqrt{\epsilon}, \mathbf{k}))$$

- 1 Introduction
- 2 Online Algorithms with ML Advice
- 3 The Predictive Marker Algorithm
- 4 Extensions
- 5 Experiments

Free parameter in the algorithm

- So far, we chose H_k as a switching point for the algorithm.

Free parameter in the algorithm

- So far, we chose H_k as a switching point for the algorithm.
- What about γH_k for some $\gamma > 0$?

Free parameter in the algorithm

- So far, we chose H_k as a switching point for the algorithm.
- What about γH_k for some $\gamma > 0$?

Theorem

Suppose that, for $\gamma > 0$, the algorithm uses γH_k as switching point. denote this algorithm by $\mathcal{PM}(\gamma)$.

Free parameter in the algorithm

- So far, we chose H_k as a switching point for the algorithm.
- What about γH_k for some $\gamma > 0$?

Theorem

Suppose that, for $\gamma > 0$, the algorithm uses γH_k as switching point. denote this algorithm by $\mathcal{PM}(\gamma)$.

Then, the competitive ratio of ϵ -assisted ϵ -assisted $\mathcal{PM}(\gamma)$ is bounded by:

$$CR_{\mathcal{PM}(\gamma), \ell}(\epsilon) \leq 2 \min \left(1 + \frac{1 + \gamma}{\gamma} \sqrt{5\epsilon}, (1 + \gamma) H_k, k \right)$$

Free parameter in the algorithm

- So far, we chose H_k as a switching point for the algorithm.
- What about γH_k for some $\gamma > 0$?

Theorem

Suppose that, for $\gamma > 0$, the algorithm uses γH_k as switching point. denote this algorithm by $\mathcal{PM}(\gamma)$.

Then, the competitive ratio of ϵ -assisted ϵ -assisted $\mathcal{PM}(\gamma)$ is bounded by:

$$CR_{\mathcal{PM}(\gamma), \ell}(\epsilon) \leq 2 \min \left(1 + \frac{1 + \gamma}{\gamma} \sqrt{5\epsilon}, (1 + \gamma) H_k, k \right)$$

- What does a low γ mean?

Free parameter in the algorithm

- So far, we chose H_k as a switching point for the algorithm.
- What about γH_k for some $\gamma > 0$?

Theorem

Suppose that, for $\gamma > 0$, the algorithm uses γH_k as switching point. denote this algorithm by $\mathcal{PM}(\gamma)$.

Then, the competitive ratio of ϵ -assisted ϵ -assisted $\mathcal{PM}(\gamma)$ is bounded by:

$$CR_{\mathcal{PM}(\gamma), \ell}(\epsilon) \leq 2 \min \left(1 + \frac{1 + \gamma}{\gamma} \sqrt{5\epsilon}, (1 + \gamma) H_k, k \right)$$

- What does a low γ mean?
- What does a high γ mean?

Robustifying LRU

- How can we use LRU in the predictive marker setting?

Robustifying LRU

- How can we use LRU in the predictive marker setting?
- Define $h(\sigma_i) = -i$.

Robustifying LRU

- How can we use LRU in the predictive marker setting?
- Define $h(\sigma_i) = -i$.

Fact

If we never switched to random evictions, this is exactly LRU.

Robustifying LRU

- How can we use LRU in the predictive marker setting?
- Define $h(\sigma_i) = -i$.

Fact

If we never switched to random evictions, this is exactly LRU.

- LRU is deterministic and therefore has only a bound of $\Theta(k)$, but is very good in practice.

Robustifying LRU

- How can we use LRU in the predictive marker setting?
- Define $h(\sigma_i) = -i$.

Fact

If we never switched to random evictions, this is exactly LRU.

- LRU is deterministic and therefore has only a bound of $\Theta(k)$, but is very good in practice.
- This new setting reduces the analysis of LRU from $\Theta(k)$ to $O(\log k)$, while still exploiting its predictive power.

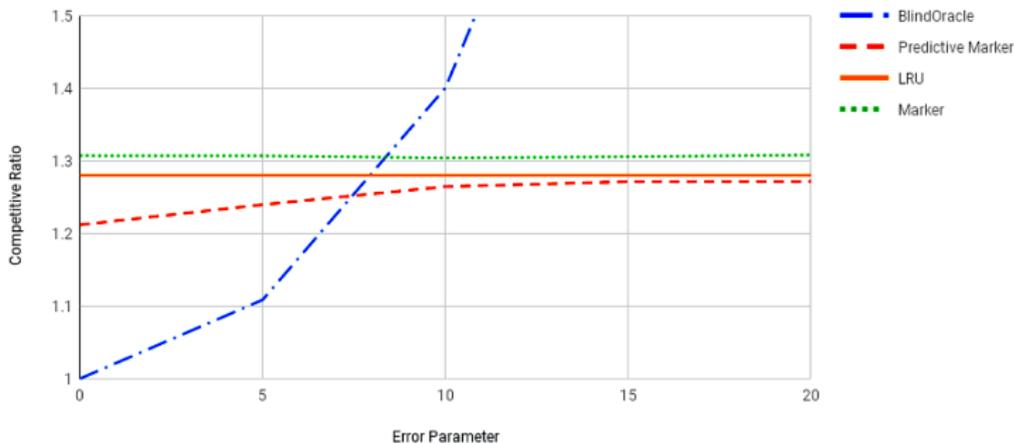
- 1 Introduction
- 2 Online Algorithms with ML Advice
- 3 The Predictive Marker Algorithm
- 4 Extensions
- 5 Experiments

Comparing results

- Lets take a look at some real world datasets:

Dataset	Num Sequences	Sequence Length	Unique Elements
BK	100	2,101	67– 800
Citi	24	25,000	593 – 719

Comparing results



Algorithm	Competitive Ratio on BK	Competitive Ratio on Citi
Blind Oracle	2.049	2.023
LRU	1.280	1.859
Marker	1.310	1.869
Predictive Marker	1.266	1.810

- We saw a general framework for combining **online algorithms** with guidance of **machine learning** - OMLA.

- We saw a general framework for combining **online algorithms** with guidance of **machine learning** - OMLA.
- We saw how we can maintain **robustness** while exploiting predictions to improve **competitiveness**.

- We saw a general framework for combining **online algorithms** with guidance of **machine learning** - OMLA.
- We saw how we can maintain **robustness** while exploiting predictions to improve **competitiveness**.
- We discussed the analysis of **robustness vs competitiveness** and looked at real-world examples.

- We saw a general framework for combining **online algorithms** with guidance of **machine learning** - OMLA.
 - We saw how we can maintain **robustness** while exploiting predictions to improve **competitiveness**.
 - We discussed the analysis of **robustness vs competitiveness** and looked at real-world examples.
-
- Thank you!



Thodoris Lykouris and Sergei Vassilvitskii.

2021. Competitive Caching with Machine Learned Advice.
J. ACM 68, 4, Article 24 (July 2021), 25 pages.