

# Online Learning For Pandora's Box and Min Sum Set Cover

## A Search Problem

Driving to work, which route to choose? (Traffic/road closures etc)



Spend time finding route with less traffic

Find the best alternative!



- Information is not free
  - Explore alternatives (open boxes)
  - Stop anytime and take best so far
- Strategy**

**Goal:**  
Find minimum cost strategy

Opening cost: 3  
Final option: Box 4  
**Total cost: 35+3**

## Our setting: repeated search problem

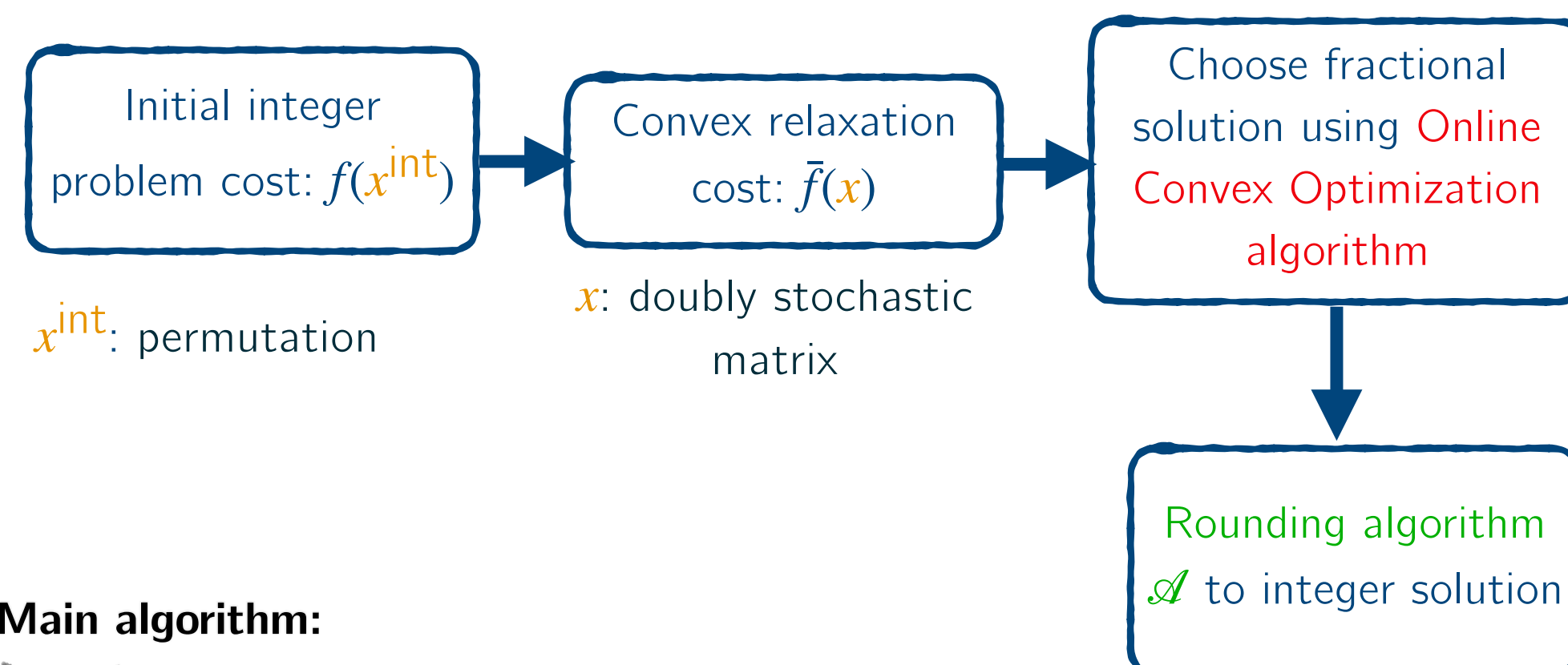
For each day  $t$ :

- New realization of values in boxes
- Pick a strategy  $\mathcal{A}^t$
- Play  $x_t$  according to  $\mathcal{A}^t$
- Receive loss  $f^t(x_t)$
- See loss function  $f^t$  on all  $x$ 's

**Goal:**  
Obtain  $\alpha$ -approximate no regret algorithm vs hindsight optimal

Full Information Setting

## Our framework



**Main algorithm:**

- $f^t(x)$  is the fractional objective function
- For each round  $t \in [T]$  do
  - Set  $x_t = \text{OCO}(f^1, \dots, f^{t-1})$
  - Round  $x_t$  to  $x_t^{\text{int}}$  according to algorithm  $\mathcal{A}$
  - Receive loss  $f^t(x_t^{\text{int}})$

Regret guarantee  $R(T)$  +  $\alpha$ -approx rounding algorithm = **Theorem 3.1**  $\alpha R(T)$  approx. regret

Optimize the two components independently!

- Use Follow The Regularized Leader family
- Set  $\text{OCO} = \min_x \sum_{\tau=1}^{t-1} f^\tau(x) + \text{Regularizer}(x)$
- Choose regularizer =  $\sum_{i,t} x_{it} \log x_{it}$
- Use randomized rounding<sup>[1]</sup>
- Guarantee  $\mathbb{E}[\bar{f}(x)] \leq \alpha f(x^{\text{int}})$
- Does not depend on  $f^t$

**Corollary 3.3.1**  
Algorithm is 9.22-approx no regret

## Bandit Setting

For each day  $t$ :

- New realization of values in boxes
- Pick a strategy  $\mathcal{A}^t$
- Play  $x_t$  according to  $\mathcal{A}^t$
- Receive loss  $f^t(x_t)$
- See loss function **only** on  $x_t$

**Idea:**  
Balance explore/FTRL steps

**Main algorithm:**

- Split  $[T]$  into intervals  $\mathcal{I}_i$ , choose uniformly random  $t_p \in [\mathcal{I}_i]$ ,  $\mathcal{R} = \emptyset$
- For each interval  $\mathcal{I}_i$  and each time  $t \in \mathcal{I}_i$ 
  - If  $t = t_p$ 
    - Open all boxes, include  $t_p$  in  $\mathcal{R}$
  - Else
    - Set  $x_t = \min_x \sum_{\tau \in \mathcal{R}} f^\tau(x) + \text{Regularizer}(x)$
    - Round  $x_t$  to  $x_t^{\text{int}}$  according to algorithm  $\mathcal{A}$

**Theorem 4.1**  
In the bandit setting, OCO Algorithm is no regret

## Summary of Results

	1 box	k boxes	Matroid basis, size k
<b>Full information &amp; bandit Against PA</b>			
$\alpha$ -approx. Regret	$\alpha = 9.22$	$\alpha = O(1)$	$\alpha = O(\log k)$
<b>Against NA</b>			
$\alpha$ -approx. Regret	$\alpha = 3.16$	$\alpha = 12.64$	$\alpha = O(\log k)$

Different ellipsoid-based algorithm