

# Finding Optimal Deterministic Policies for **Constrained Stochastic Shortest Path Problems**



## Johannes Schmalz and Felipe Trevizan

### Background

### Constrained Stochastic Shortest Path Problem (CSSP)

Defined by tuple  $(S, s_0, G, A, P, \vec{C}, \vec{u})$  with states (S), initial state  $(s_0)$ , goals (G), actions (A), transition probabilities (P), cost vector  $(\vec{C}: \mathsf{S} \to \mathbb{R}^{n+1}_{>0})$ , cost upper bounds  $(\vec{u} \in \mathbb{R}^{n}_{>0})$ .

### Policies

Policies  $\pi$  map states to actions. They have two flavours:

**Deterministic** policy  $\pi : S \rightarrow A$ 

**Stochastic** policy  $\pi : S \rightarrow$  probability distribution over A

**Policy Cost**:  $C_i(\pi) =$ expected cost *i* to reach G with  $\pi$ 

### **CSSP** Example

You need to get to work. You can run, use a taxi, or walk to the train station and try the train. The train is cancelled with 50% probability. Each action has a cost vector  $[t \ p \ e]$  in terms of time (t), price (p), and personal effort (e). Task: get to work in minimal time s.t. price  $\leq 15$  and effort  $\leq 10$  over expectation.

### run [1 0 20]



### Solving CSSPs with MIP

Imagine the CSSP as a flow network: actions are pipes, states are junctions, and we want to route 1L of water through it with minimal cost. Each  $x_{s,a}$  denotes amount of flow through pipe afrom s and  $\Delta_{s,a} = 1/0$  if the pipe has/has no flow.

 $\min C_0(\vec{x})$  s.t.  $\vec{x}.\vec{\Lambda}$  $\forall s \in \mathsf{S} \setminus (\mathsf{G} \cup \{s_0\})$ out(s) - in(s) = 0 $out(s_0) - in(s_0) = 1$  $\sum in(g) = 1$ 

**Feasibility**:  $\pi$  must satisfy  $C_i(\pi) \leq u_i$  for all  $i \in \{1, \ldots, n\}$ 

**Optimality**:  $\pi$  is optimal if it is feasible and minimises  $C_0(\pi)$ 

Optimal stochastic policies are optimal for the CSSP, but optimal deterministic policies need **not** be optimal for the CSSP. Nevertheless...

### **Practitioners Want Deterministic Policies**

- ethical issues in medical contexts (Roijers et al. 2013)
- aviation regulations (Geißer et al. 2020)
- coordination in multi-agent systems (Dolgov and Durfee) 2005)
- accountability and explainability (Hong and Williams 2023)

more predictable (see example)

**Optimal stochastic policy:**  $\pi^*(s_0) = \{ \operatorname{run} : 50\%, \operatorname{taxi} : 50\% \}$  $C_{\rm price}(\pi^*) = 15$   $C_{\rm effort}(\pi^*) = 10$  $C_{\mathsf{time}}(\pi^*) = 1$ But individually *run* and *taxi* violate the constraints...

**Optimal deterministic policy:**  $\pi(s) =$  walk or train  $C_{\mathsf{price}}(\pi) = 5$  $C_{\mathsf{time}}(\pi) = 3$  $C_{\text{effort}}(\pi) = 4$ More expensive time-wise, but satisfies constraints \*

$g\inG$	
$x_{s,a} \ge 0$	$\forall s \in S, a \in A(s)$
$C_i(\vec{x}) \le u_i$	$\forall i \in \{1, \dots, n\}$
$x_{s,a} \leq \mathbf{M}\Delta_{s,a}$	$\forall s \in S, a \in A(s)$
$\sum \Delta_{s,a} \le 1$	$\forall s \in S$
$a{\in}A(s)$	
$\Delta_{s,a} \in \{0,1\}$	$\forall s \in S, a \in A(s)$

In and out flow from states, and flow cost are macros:

• 
$$C_i(\vec{x}) = \sum_{s \in S, a \in A(s)} x_{s,a} C_i(a)$$
  
•  $out(s) = \sum_{a \in A(s)} x_{s,a}$   
•  $in(s) = \sum_{s' \in S, a' \in A(s')} x_{s,a} P(s'|s,a)$ 

Contributions

### The Issue with big $\mathbf{M}$

If M is too small...

■ MIP can become infeasible

optimal policy may become infeasible If M is too big...

Numerical instability

May get non-integer solutions (trickle flow)

### Finding big M Automatically

**New insight:** we can relate **any feasible** solution  $\vec{x}$  to the maximum flow over the **optimal** solution with  $obj(\vec{x}) \cdot g^{-1} \ge 1$  $x_{\max}$  where  $g = \min_{a \in A} C_0(a)$ . Algorithm: 1 select some M **2** try to solve MIP with  $\mathbf{M}$ 

### New Algorithm for Finding Deterministic Policies

For finding **stochastic** policies, i<sup>2</sup>-dual (Trevizan et al. 2017) is the state-of-the-art. It uses a heuristic to iteratively construct partial CSSPs, focusing on the promising states. The partial CSSPs are solved with LPs.

### A General Bound

This always works:  $\mathbf{M} = p_{\min}^{-|\mathsf{S} \setminus \mathsf{G}|}$ ■  $p_{\min} \in (0, 1]$  min. probability in CSSP **\blacksquare** S \ G set of non-goal states

That's impractically big, but can't do better in general:



 $\blacksquare$  if infeasible: increase  ${\bf M}$  and repeat step 2 • if feasible: set  $\mathbf{M} \leftarrow \operatorname{obj}(\vec{x}) \cdot g^{-1}$  ${}_{\mathbf{3}}$  solve MIP with  ${\mathbf{M}}$ 

### Avoiding big M

Can completely avoid big  ${f M}$  with SOS1 constraints. A SOS1 constraint is an ordered set of continuous variables  $\{x_0, \ldots, x_k\}$ such that at most one variable is allowed to be nonzero.

For us:  $\{x_{s,a} | a \in A(s)\}$  for each state s, i.e., at most one action for A(s) may have nonzero flow. That's exactly what we want!

**New:** we replace  $i^2$ -dual's LPs with MIPs. This yields i<sup>2</sup>-dual-det, which finds optimal **deterministic** policies.

Making it faster: we don't care about the exact solution to each MIP, so we can approximate them! There are many ways to do this:

- Use LP relaxation
- Use constraint generation for integrality constraints
- Allow large MIP gap (also makes it anytime)

### **Interesting Benchmarks and Performance**

### More Background: Linearisation

A linearisation  $\vec{\lambda} \in \mathbb{R}^{n+1}$  relaxes the CSSP into an SSP with the scalar cost function  $C'(a) = \lambda \cdot \hat{C}(a)$ 

 $\dot{\lambda} = [1 \ 0 \ \cdots \ 0]$  is the trivial relaxation that only looks at  $C_0$ .

### **Constraint Interestingess**

We introduce a **new categorisation of CSSPs** based on how interesting their constraints are.



- **Trivial**: if the constraints can be ignored, and the trivial relaxation yields a feasible policy.
- **Linearisable**: if  $\exists \hat{\lambda}$  so that its relaxation yields a feasible policy.
- **Interesting**: if it is not linearisable.
- Practically Interesting: if it is linearisable, but for each feasible policy there are many infeasible ones.

Policy costs w.r.t.  $C_0$  and  $C_1$ . Red line is constraint over  $C_1$ . Dotted green line is contour line for linearisation. Points are the costs of deterministic policies.  $\blacklozenge / \circ =$  feasible / infeasible w.r.t.  $C_2$ 

### **Very Short Summary of Experimental Results**

• The existing state-of-the-art, which is based on linearisation, is best on **trivial** and **linearisable** problems. Our algorithm i<sup>2</sup>-dual-det is best on interesting problems.

More available at schmlz.github.io/det-pi-for-cssp

You can email me at johannes.schmalz@anu.edu.au

