

# Analyse eines Forschungsthemas

## Stochastic Shortest Paths

Maximilian Starke

Fakultät für Informatik  
Technische Universität Dresden

5. September 2022

Introduction

Essential Definitions

The classic stochastic shortest path problem

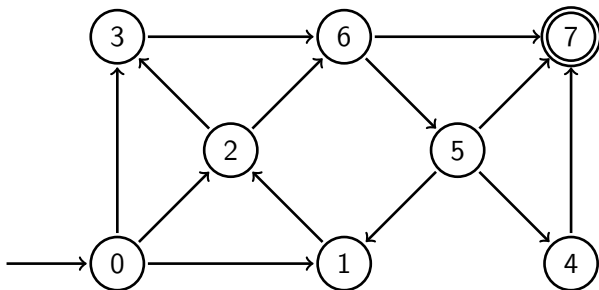
Different variants of the stochastic shortest path problem

Keep an eye on the variance

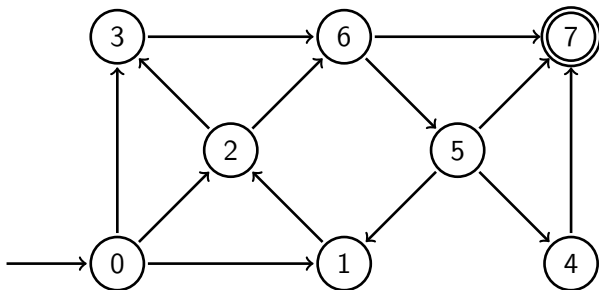
# Section 1

## Introduction

- The *simplest* shortest path problem

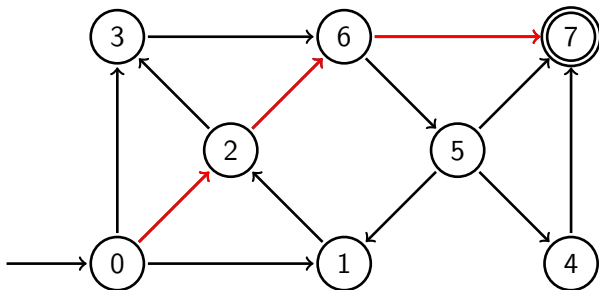


- The *simplest* shortest path problem



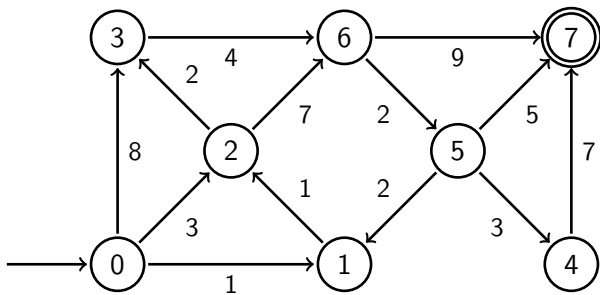
- Task
  - Find the shortest path (*number of hops*)!

- The *simplest* shortest path problem

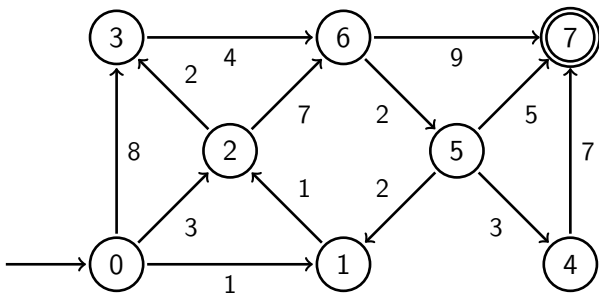


- Task
  - Find the shortest path (*number of hops*)!

- The *classical, non-stochastic, deterministic* shortest path problem



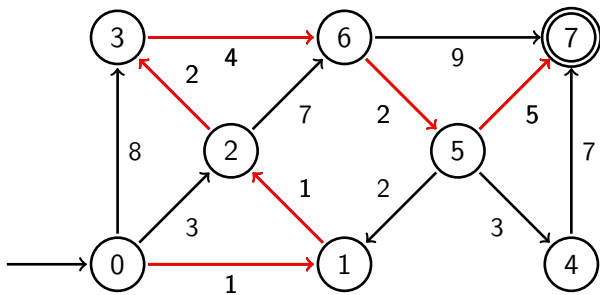
- The *classical, non-stochastic, deterministic* shortest path problem



- Task
  - Find the path with the minimal weight sum!

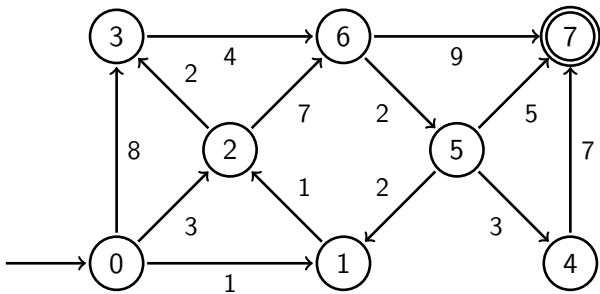


- The *classical, non-stochastic, deterministic* shortest path problem



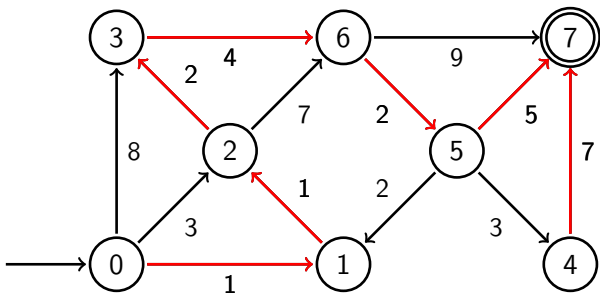
- Task
  - Find the path with the minimal weight sum!

- The *classical, non-stochastic, deterministic* shortest path problem



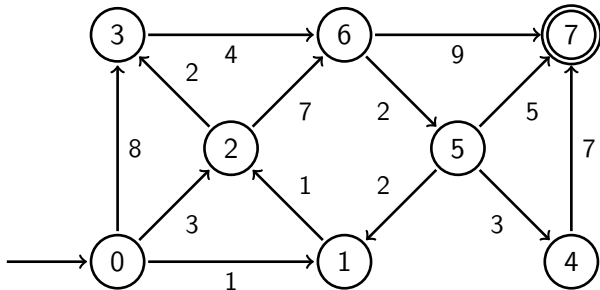
- Task
  - Find the path with the minimal weight sum!
  - Give a strategy to always reach the goal while collecting minimal weight!

- The *classical, non-stochastic, deterministic* shortest path problem

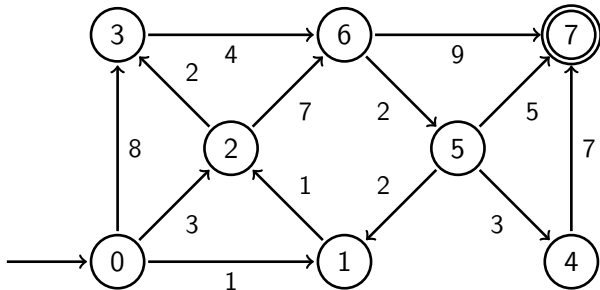


- Task
  - Find the path with the minimal weight sum!
  - Give a strategy to always reach the goal while collecting minimal weight!

- The *stochastic* shortest path problem

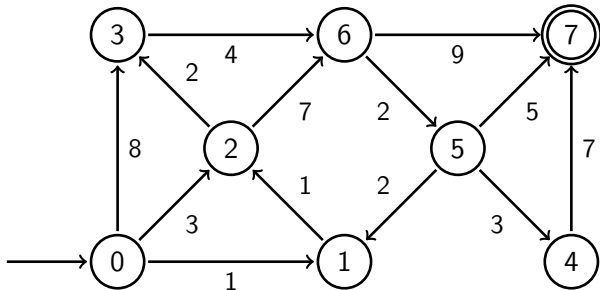


- The *stochastic* shortest path problem



- Markov Decision Process (MDP)

- ▶ The *stochastic* shortest path problem



- ▶ Markov Decision Process (MDP)
- ▶ Task
  - ▶ Give a strategy to reach the goal with minimal *expected* accumulated weights!

## Section 2

### Essential Definitions

- ▶ MDP
- ▶ Expectation
- ▶ Conditional Expectation
- ▶ Variance-penalized Expectation
- ▶ schedulers, kind of schedulers...



# The classic stochastic shortest path problem

- ▶ given:
  - ▶ a single goal state
  - ▶ positive cycle condition: There is no cycle with  $\sum wgt \geq 0$
  - ▶ goal is reachable from each state
- ▶ goal: Maximize the expected accumulated weight until reaching goal state.
- ▶ Well known for a long time:
  - ▶ There exists an optimal memoryless deterministic scheduler  $\mathfrak{S}$ .
  - ▶  $\mathfrak{S}$  is computable by solving a LP
  - ▶ iterative algorithm:
    - ▶ start at any feasible scheduler
    - ▶ iterative improvement
    - ▶ stop at an optimal vertex of the LP (corresponding to some MD scheduler)

## Section 3

### The classic stochastic shortest path problem

# The classic stochastic shortest path problem

- ▶ Can we do it better?
- ▶ YES! - using **spider construction**!
- ▶ Assume furthermore:
  - ▶  $\mathcal{M}$  is an MDP with arbitrary integer weights
- ▶ The following can be solved in polynomial time:
  - ▶ Check:  $\mathbb{E}_{\mathcal{M},s}^{inf}(\boxplus \text{goal}) > -\infty$ ?
  - ▶ Compute  $\mathbb{E}_{\mathcal{M},s}^{inf}$  if it is finite

# Spider Construction

- ▶ Idea: construct a new MDP  $\mathcal{N}$  from the given MDP  $\mathcal{M}$
- ▶ Pick a 0-BSCC  $\mathcal{E}$  of  $\mathcal{M}$  and some vertex  $s_0$  in  $\mathcal{E}$ .
- ▶  $\mathcal{M} \mapsto \mathcal{N} := \text{Spider}_{\mathcal{E}, s_0}(\mathcal{M})$
- ▶ The spider construction is done by applying the following steps:
  1. Remove all actions  $(s, s_\alpha) \in \mathcal{E}$
  2. Add actions  $(s, \tau)$  for all  $s \in \mathcal{E} \setminus \{s_0\}$  such that
    - ▶  $P_{\mathcal{N}}(s, \tau, s_0) := 1$
    - ▶  $\text{wgt}_{\mathcal{N}}(s, \tau) := \text{wgt}_{\mathcal{M}}(s, s_0)$
  3. For each  $s \in \mathcal{E} \setminus \{s_0\}$  and  $\beta \in \text{Act}_{\mathcal{M}}(s) \setminus \{\alpha_s\}$  let us replace  $(s, \beta)$  by  $(s_0, \beta)$  where
    - ▶  $P_{\mathcal{N}}(s_0, \beta, u) := P_{\mathcal{M}}(s, \beta, u)$
    - ▶  $\text{wgt}_{\mathcal{N}}(s_0, \beta) + \text{wgt}_{\mathcal{N}}(s, \tau) = \text{wgt}_{\mathcal{M}}(s, \beta)$

# Classification of paths

A path  $\pi \in \text{InfPaths}(\mathcal{M})$  is called

- ▶ pumping  $:\Leftrightarrow \liminf_{n \rightarrow \infty} (\text{wgt}(\text{pref}(\pi, n))) = \infty$
- ▶ (positively)  
negatively weight divergent  $:\Leftrightarrow \begin{array}{l} \limsup_{n \rightarrow \infty} \\ \liminf_{n \rightarrow \infty} \end{array} = \begin{array}{l} \infty \\ -\infty \end{array}$
- ▶ gambling  $:\Leftrightarrow \pi$  is positively and negatively weight divergent
- ▶ bounded from below  $:\Leftrightarrow \liminf_{n \rightarrow \infty} \text{wgt}(\text{pref}(\pi, n)) \in \mathbb{Z}$

# Classification of end components

We distinguish end components by the following types

- ▶ pumping ECs:  $\exists$  scheduler  $\mathcal{G} : \Pr(\pi \text{ is pumping}) = 1$
- ▶ (positively)  
negatively weight divergent ECs:  $\exists$  scheduler  $\mathcal{G} : \Pr(\pi \text{ is (positively)  
negatively weight divergent}) = 1$
- ▶ gambling ECs:  $\mathbb{E}(\text{MP}) = 0$  and it is positively and negatively weight divergent
- ▶ bounded EC: There exists an upper bound and a lower bound

## Section 4

### Different variants of the stochastic shortest path problem

Imagine to not reach goal with probability 1.

- ▶ maximal conditional expected accumulated reward
- ▶ partial expected accumulated reward.



conditional expectation	partial expectation
$\mathbb{CE}$	$\mathbb{PE}$
$\mathbb{CE} = \mathbb{E}(\boxplus \text{goal} \mid \Diamond \text{goal})$	$\pi \not\models \Diamond \text{goal} \Rightarrow \text{wgt}(\pi) := 0$
good approximation for maximizing probability and reward until goal	may lead to quite high $\mathbb{CE}$ paired with low probability of reaching goal

# conditional expected accumulated reward

Given:

- ▶ MDP  $\mathcal{M}$  with non-negative integer weights
- ▶ two sets of states  $F, G \subseteq \text{States}(\mathcal{M})$

## Definition

$$\mathbb{CE}^{max} := \sup_{\mathfrak{G} \in S} (\mathcal{E}_{\mathcal{M}, s_i \text{ init}}^{\mathfrak{G}}(\boxplus F \mid \Diamond G))$$

where  $S$  is the set of schedulers:

$$S := \{\mathfrak{G} \mid \Pr_{\mathcal{M}, s_i \text{ init}}^{\mathfrak{G}}(\Diamond G) > 0 \wedge \Pr_{\mathcal{M}, s_i \text{ init}}^{\mathfrak{G}}(\Diamond F \mid \Diamond G) = 1\}$$

# results about conditional expected accumulated rewards

- ▶ There is a **PTime algorithm** to decide: Is  $\mathbb{CE}^{max}$  finite?
- ▶ There is a **pseudo-PTime algorithm** to calculate an upperbound  $\mathbb{CE}^{ub} \geq \mathbb{CE}^{max}$

- ▶ If we have  $F = G$  and

$\forall s \in \text{States}(\mathcal{M}) : s \models \exists \Diamond G \Rightarrow \Pr_{\mathcal{M},s}^{min}(\Diamond G) > 0$  there is a

**PTime algorithm** to calculate an upperbound  $\mathbb{CE}^{ub} \geq \mathbb{CE}^{max}$

- ▶ The problem Decide if  $\mathbb{CE}^{max} \bowtie t$  where we have
  - ▶  $t \in \mathbb{Q} \dots$  some rational threshold
  - ▶  $\bowtie \in \{<, \leq, \geq, >\}$

is **PSPACE-hard**, solvable in **ExpTime** and *for acyclic MDPs* **PSPACE-complete**

- ▶ In **ExpTime** we can compute  $\mathbb{CE}^{max}$  together with an optimal scheduler

## Section 5

Keep an eye on the variance