# Introduction to
# Artificial Intelligence

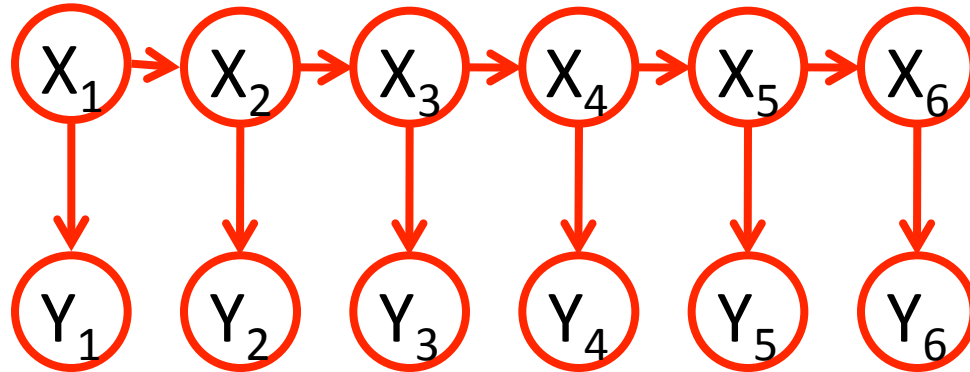## Lecture 16 – Markov Decision Processes

CS/CNS/EE 154

Andreas Krause

# Announcements

- Homework 3 out, due Wed Nov 24
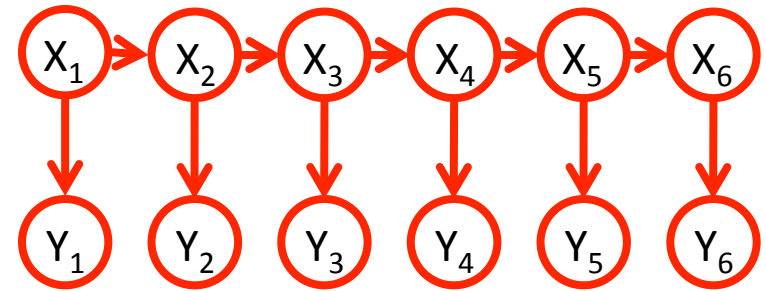- Code for project final released; due Dec 1

# HMMs / Kalman Filters



- $X_1, ..., X_T$: Unobserved (hidden) variables
- $Y_1, ..., Y_T$: Observations
- HMMs: $X_i$ Multinomial, $Y_i$ multinomial (or arbitrary)
- Kalman Filters: $X_i$, $Y_i$ Gaussian distributions

# Bayesian filtering

- Start with $P(X_1)$

- At time t

  - Assume we have $P(X_t \mid y_{1...t-1})$

  - Conditioning: $P(X_t \mid y_{1...t})$

$$P(X_t \mid y_{1:t}) = \frac{1}{Z} P(X_t \mid y_{1:t-1}) \cdot P(y_t \mid X_t)$$

Have $P(X_i)$
want $P(X_i \mid Y_i)$
$= \frac{1}{Z} P(Y_i \mid X_i) \cdot P(X_i)$
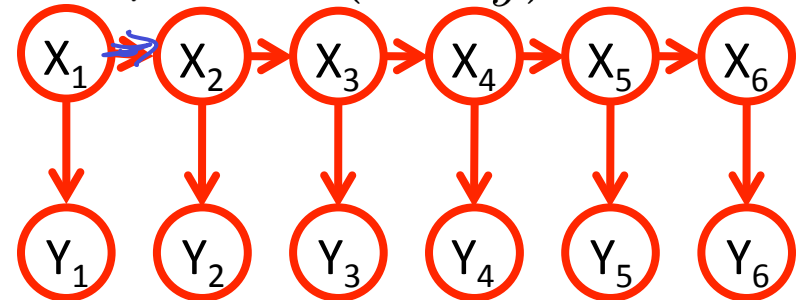Bayes' rule

  - Prediction: $P(X_{t+1} \mid y_{1...t})$

$$P(X_{t+1} \mid y_{1:t}) = \sum_{X_t} P(X_{t+1}, X_t \mid y_{1:t})$$

$$= \sum_{X_t} P(X_t \mid y_{1:t}) \underbrace{P(X_{t+1} \mid X_t, y_{1:t})}_{P(X_{t+1} \mid X_t)}$$

For k states, can do filtering in $O(t^2)$

4

# Kalman Filters (Gaussian HMMs)

- $X_1,...,X_T$: Location of object being tracked $\in \mathbb{R}^d$
- $Y_1,...,Y_T$: Observations $\in \mathbb{R}^{d'}$
- $P(X_1)$: Prior belief about location at time 1
- $P(X_{t+1}|X_t)$: "Motion model"
  - How do I expect my target to move in the environment?
  $$\mathbf{X}_{t+1} = \mathbf{F}\mathbf{X}_t + \varepsilon_t \text{ where } \varepsilon_t \in \mathcal{N}(0, \Sigma_x)$$

- $P(Y_t | X_t)$: "Sensor model"
  - What do I observe if target is at location $X_t$?
  $$\mathbf{Y}_t = \mathbf{H}\mathbf{X}_t + \eta_t \text{ where } \eta_t \in \mathcal{N}(0, \Sigma_y)$$

# General Kalman update

- Transition model $\quad P(\mathbf{x}_{t+1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t+1}; \mathbf{F}\mathbf{x}_t, \Sigma_x)$
- Sensor model $\quad\quad P(\mathbf{y}_t \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{y}_t; \mathbf{H}\mathbf{x}_t, \Sigma_y)$
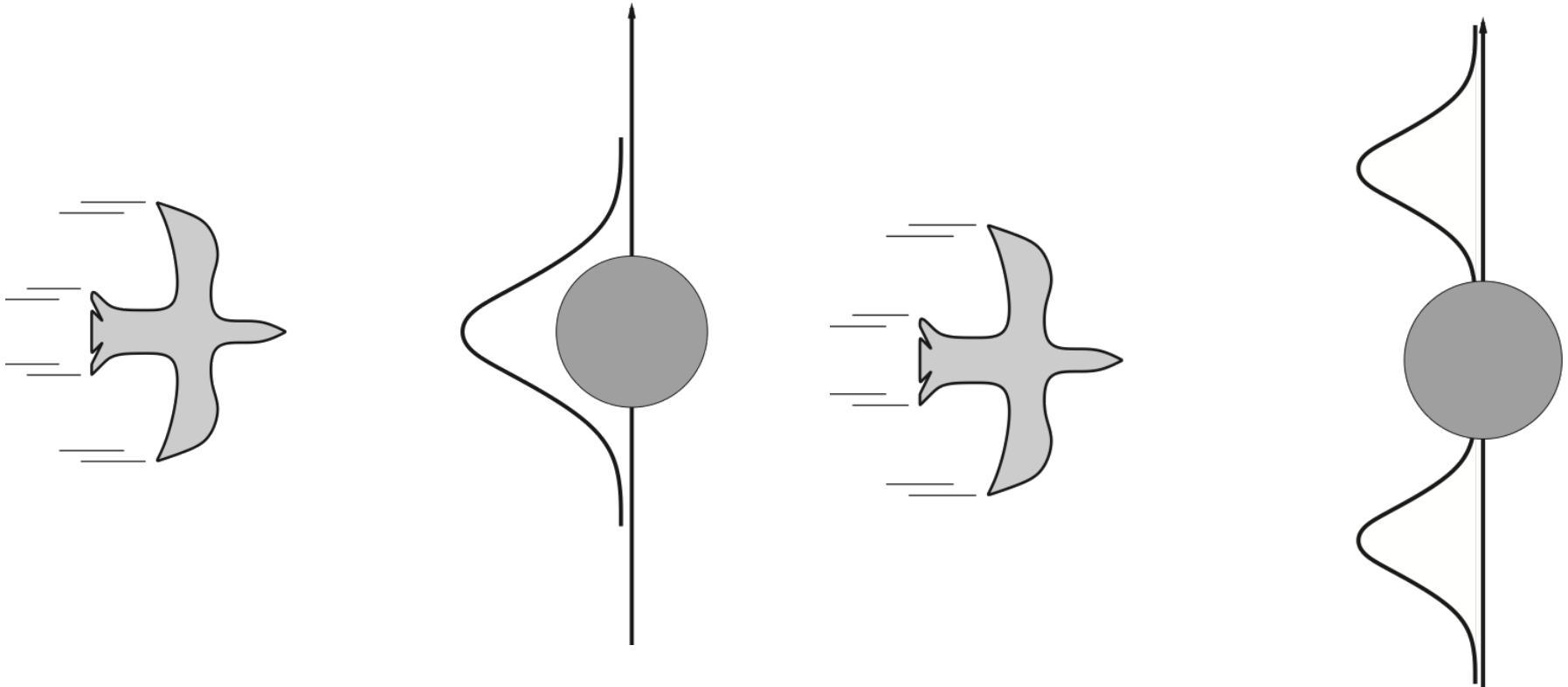
- Kalman Update: $\quad \mu_{t+1} = \mathbf{F}\mu_t + \mathbf{K}_{t+1}(\mathbf{y}_{t+1} - \mathbf{H}\mathbf{F}\mu_t)$

$$\Sigma_{t+1} = (\mathbf{I} - \mathbf{K}_{t+1})(\mathbf{F}\Sigma_t\mathbf{F}^T + \Sigma_x)$$

- Kalman gain:

$$\mathbf{K}_{t+1} = (\mathbf{F}\Sigma_t\mathbf{F}^T + \Sigma_x)\mathbf{H}^T(\mathbf{H}(\mathbf{F}\Sigma_t\mathbf{F}^T + \Sigma_x)\mathbf{H}^T + \Sigma_y)^{-1}$$

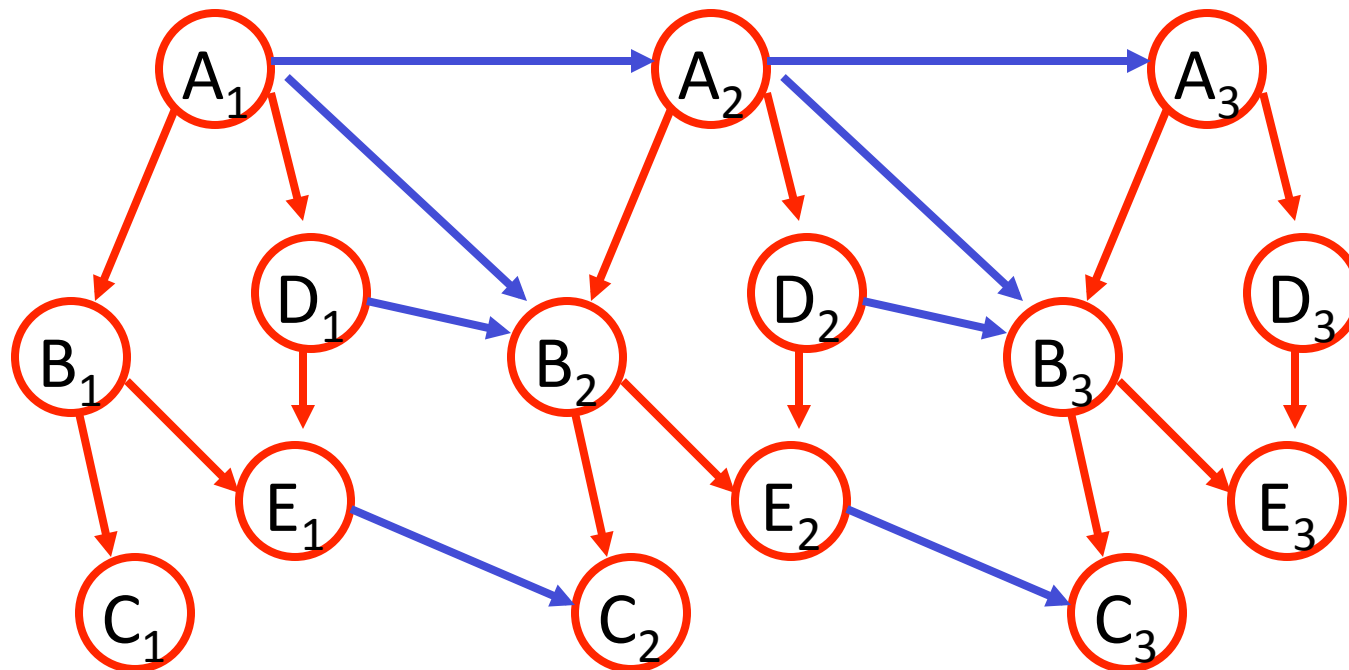- Can compute $\Sigma_t$ and $\mathbf{K}_t$ offline

# When KFs fail



- KFs assume transition model is **linear**
  - Implies that predictive distribution is Gaussian (unimodal)
- Need approximate inference to capture nonlinearities!

# Dynamic Bayesian Networks
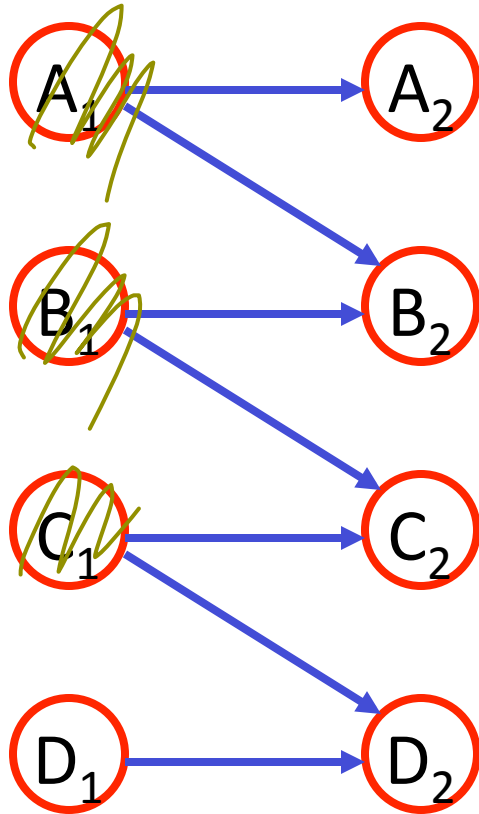
- At every timestep have a *Bayesian Network*



- Variables at each time step t called a slice $S_t$
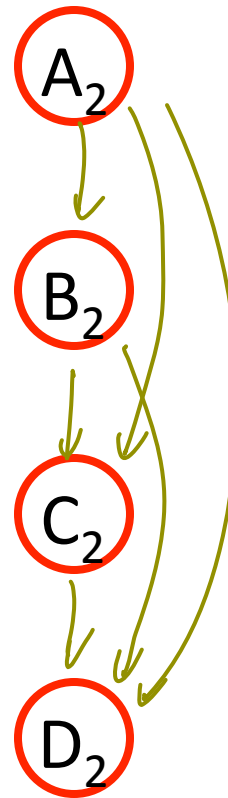- "Temporal" edges connecting $S_{t+1}$ with $S_t$

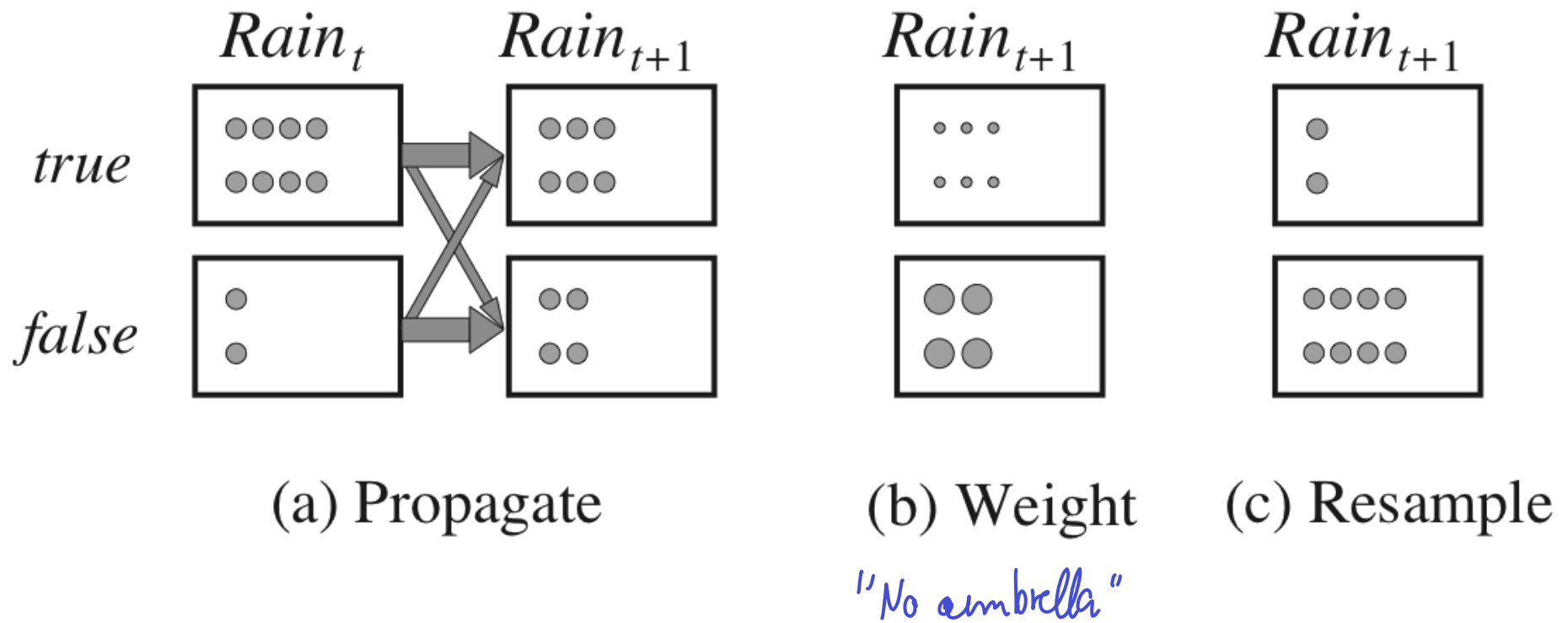# Inference in DBNs?

## DBN



## Marginals at time 2

Need approximate inference!

# Particle filtering

- Very useful approximate inference technique for dynamical models
  - Nonlinear Kalman filters
  - Dynamic Bayesian networks

- **Basic idea**: Approximate the posterior at each time by samples (particles), which are propagated and reweighted over time

# Particle filtering example



(a) Propagate

(b) Weight

(c) Resample

"No umbrella"

# Representing distributions by particles

- True distribution (possibly continuous): P(x)

- N i.i.d. samples: $x_1, \ldots, x_N$

$$\delta_{x_i}(k) = \begin{cases} 1 & \text{if } x = x_i \\ 0 & \text{otw} \end{cases}$$

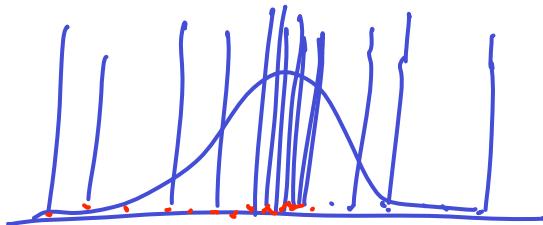- Represent: $P(x) \approx \dfrac{1}{N} \delta_{x_i}(x)$

- Get expectations: $\mathbb{E}_P[f(X)] \approx \dfrac{1}{N} \sum_i f(x_i)$

- E.g., mean: $\mathbb{E}_P[X] \approx \dfrac{1}{N} \sum_i x_i$

# Particle filtering

- Suppose

$$P(X_t \mid y_{1:t}) \approx \frac{1}{N} \sum_{i=1}^{N} \delta_{x_{i,t}}$$

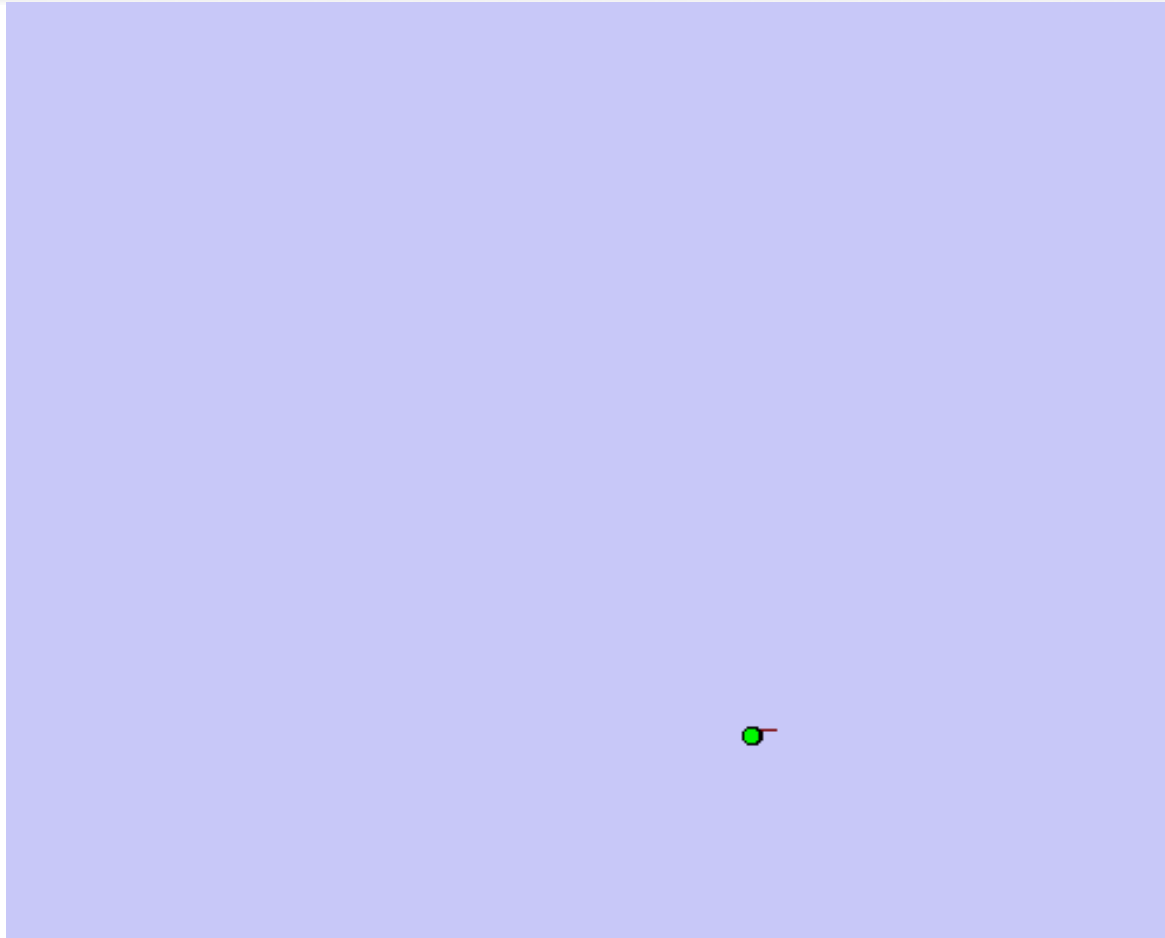- For each particle:

$$x_i' \sim P(X_{t+1} \mid x_{i,t})$$

$$Z = \sum_{i=1}^{N} P(y_{t+1} \mid x_i')$$

- Weigh particles:

$$w_i = \frac{1}{Z} P(y_{t+1} \mid x_i')$$

- Resample N particles

$$x_{i,t+1} \sim \frac{1}{N} \sum_{i=1}^{N} w_i \delta_{x_i'}$$
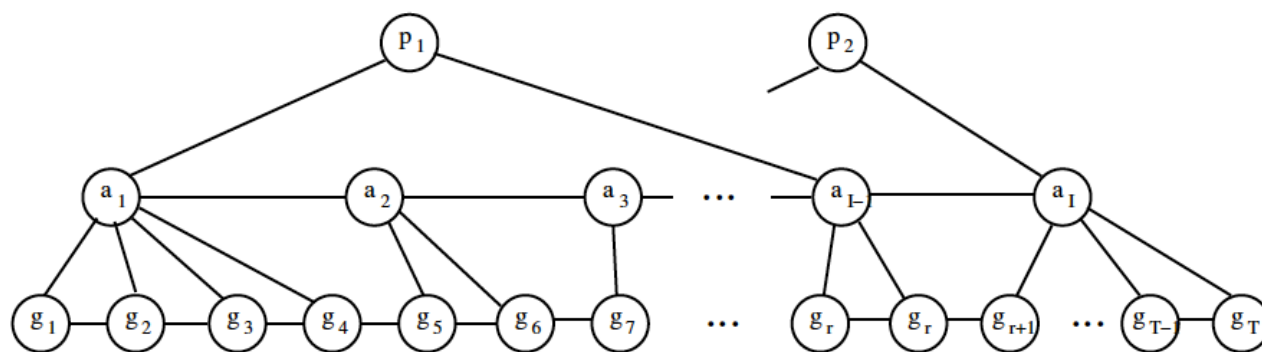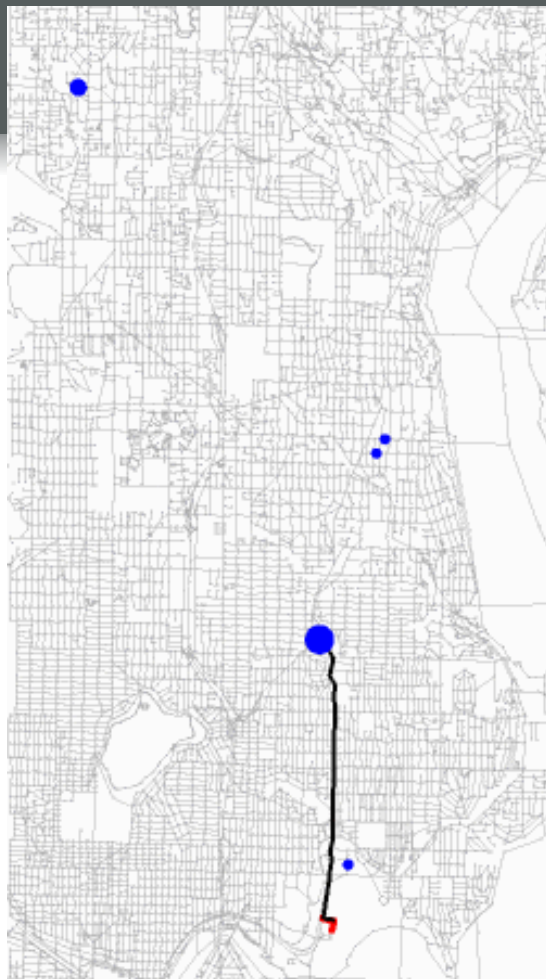
# Robot localization & mapping

D. Haehnel,
W. Burgard,
D. Fox, and
S. Thrun.
*IROS-03*.

- Infer both location and map from noisy sensor data
- Particle filters

# Activity recognition

L. Liao, D. Fox, and H. Kautz. *AAAI-04*

Predict "goals" from raw GPS data
"Hierarchical Dynamical
Bayesian networks"



**Significant places**
home, work, bus stop, parking lot, friend

**Activity sequence**
walk, drive, visit, sleep, pickup, get on bus

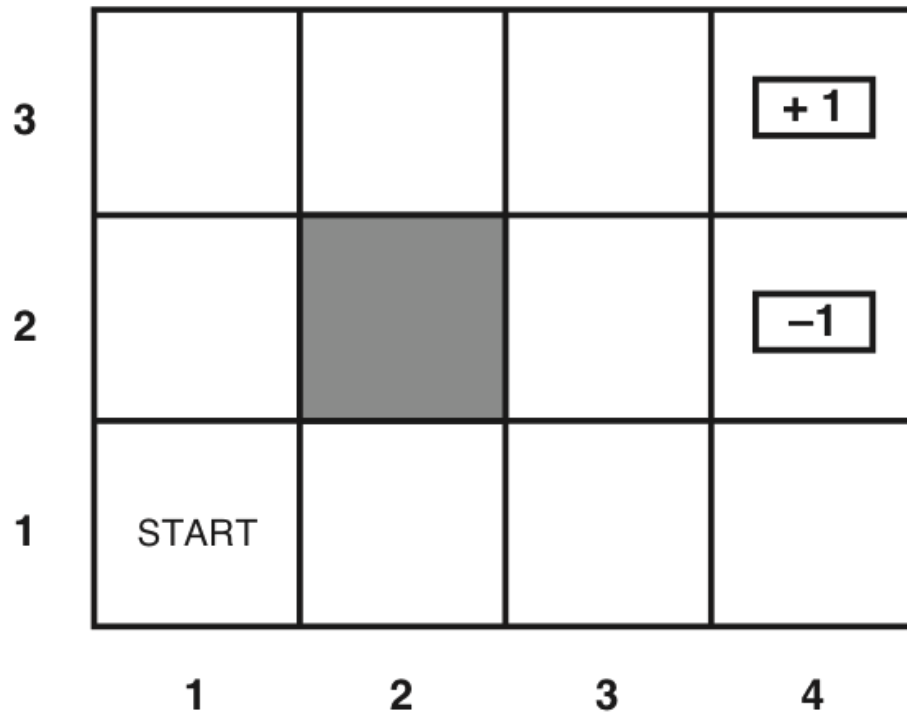**GPS trace**
association to street map

# Summary

- Dynamical models
  - Multiple copies of static models, one per time step
- Examples:
  - HMM
  - Kalman Filter
  - Dynamic Bayesian networks
- Inference tasks
  - Filtering/prediction: Can do recursively!
  - Smoothing
  - MPE
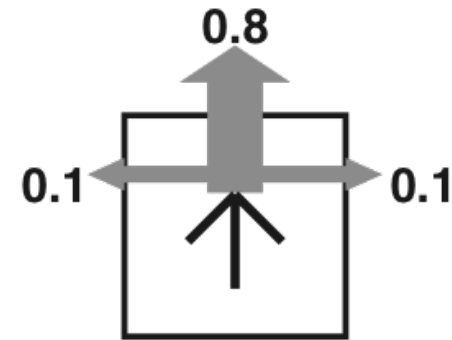- Particle filtering for approximate inference

# Probabilistic planning

- So far: Probabilistic inference in dynamical models
  - E.g.: Tracking a robot based on noisy measurements

- Next: How should we control the robot to maximize reward?
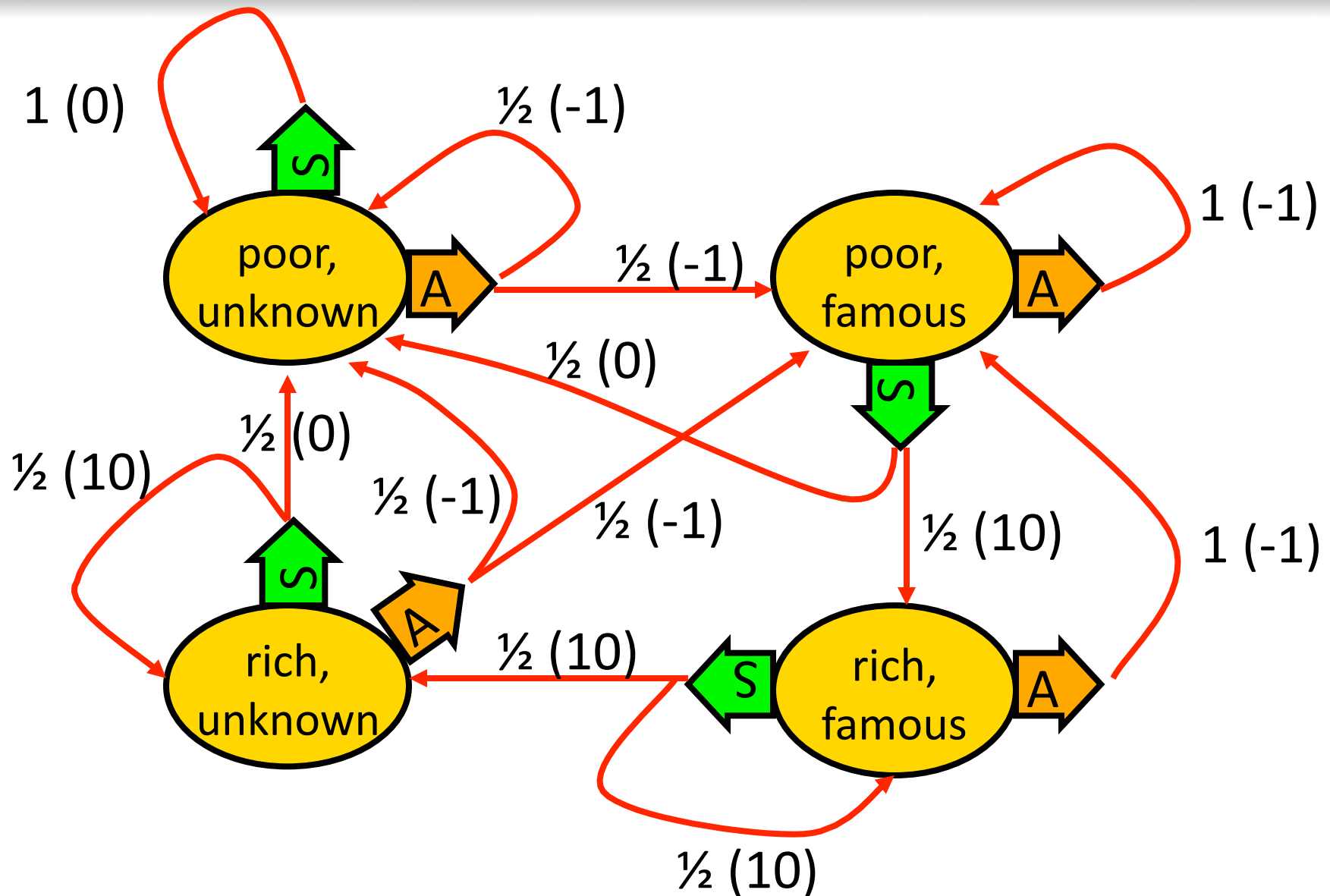
(a)
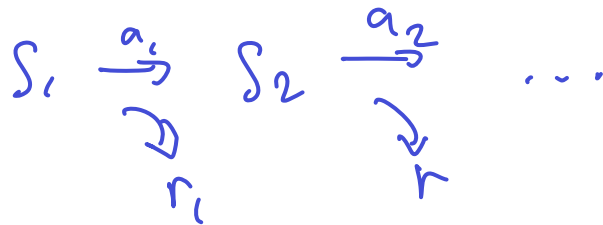
(b)

# Becoming rich and famous

# Markov Decision Processes

- An MDP has
  - A set of states $X = \{x_1, \ldots, x_n\}$ …
  - A set of actions $A = \{a_1, \ldots, a_m\}$
  - A reward function $r(x,a)$   [or random var. with mean $r(x,a)$]
  - Transition probabilities
    $P(x'|x,a) = \text{Prob}(\text{Next state} = x' \mid \text{Action } a \text{ in state } x)$

- For now assume r and P are known!

- Want to choose actions to maximize reward

# Utility over time

- Finite horizon

$$S_1 \xrightarrow{a_1} S_2 \xrightarrow{a_2} S_3 \rightarrow \cdots \rightarrow S_T$$

$$r_1 \qquad r_2 \qquad r_3 \qquad \cdots$$

$$R_T = \sum_{t=1}^{T} r_t$$

- Discounted rewards

$$S_1 \xrightarrow{a_1} S_2 \xrightarrow{a_2} \cdots$$

$$r_1 \qquad r$$

$$R_T = \sum_{t=0}^{\infty} \gamma^t r_t$$

$$\gamma \in (0,1)$$

In practice: $\gamma = 0.95$

# Finite horizon MDP Decision model

- Reward R = 0

- Start in state x

- For t = 0 to T

  - Choose action a

  - Obtain reward R = R + r(x,a)

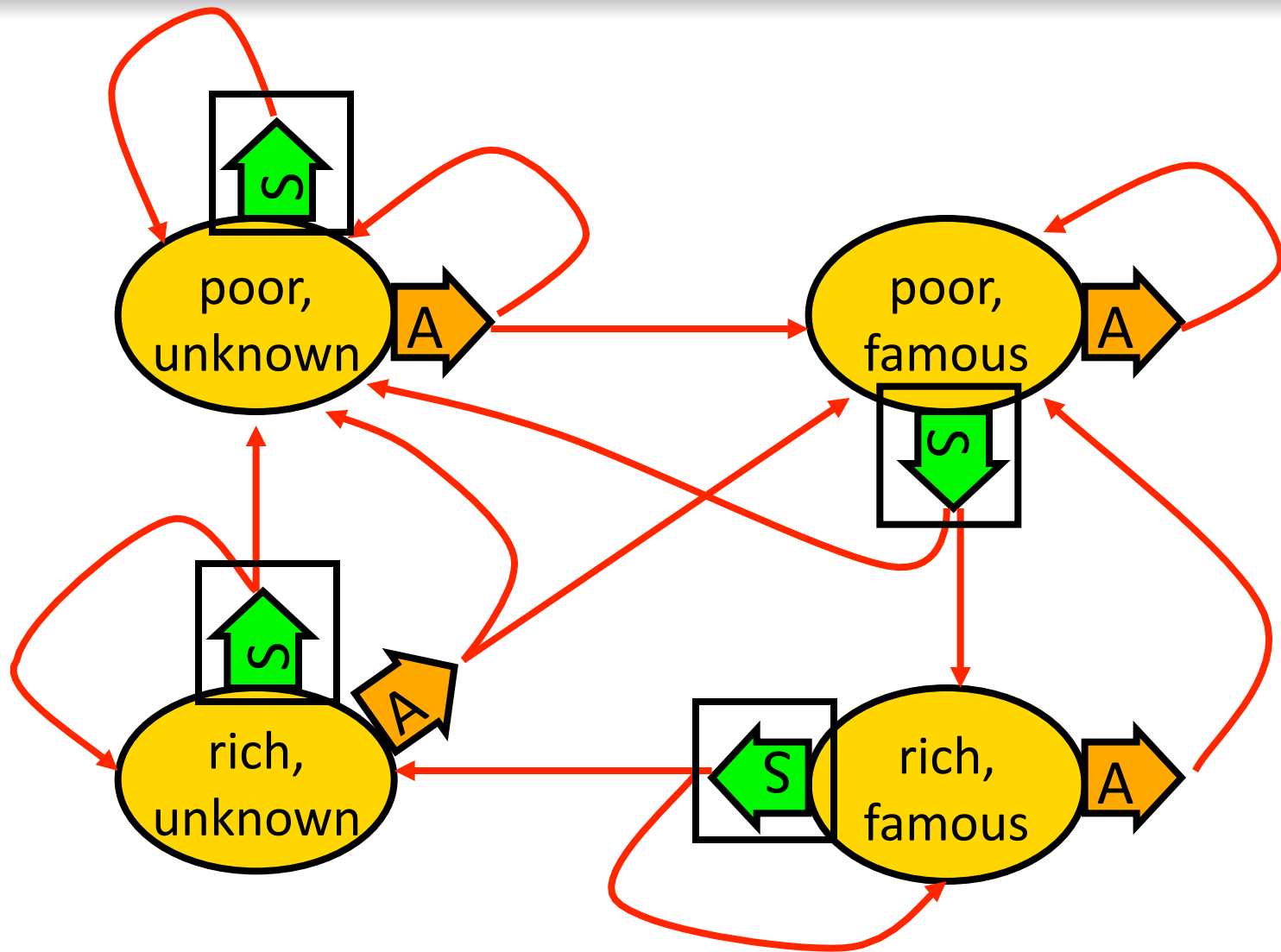  - End up in state x' according to P(x'|x,a)

  - Repeat with x ← x'

# Discounted MDP Decision model

- Reward R = 0

- Start in state x

- For t = 0 to $\boxed{\infty}$
  - Choose action a
  - Obtain **discounted** reward R = R + $\boxed{\gamma^t}$ r(x,a)
  - End up in state x' according to P(x'|x,a)
  - Repeat with x ← x'

**This lecture**: Discounted rewards
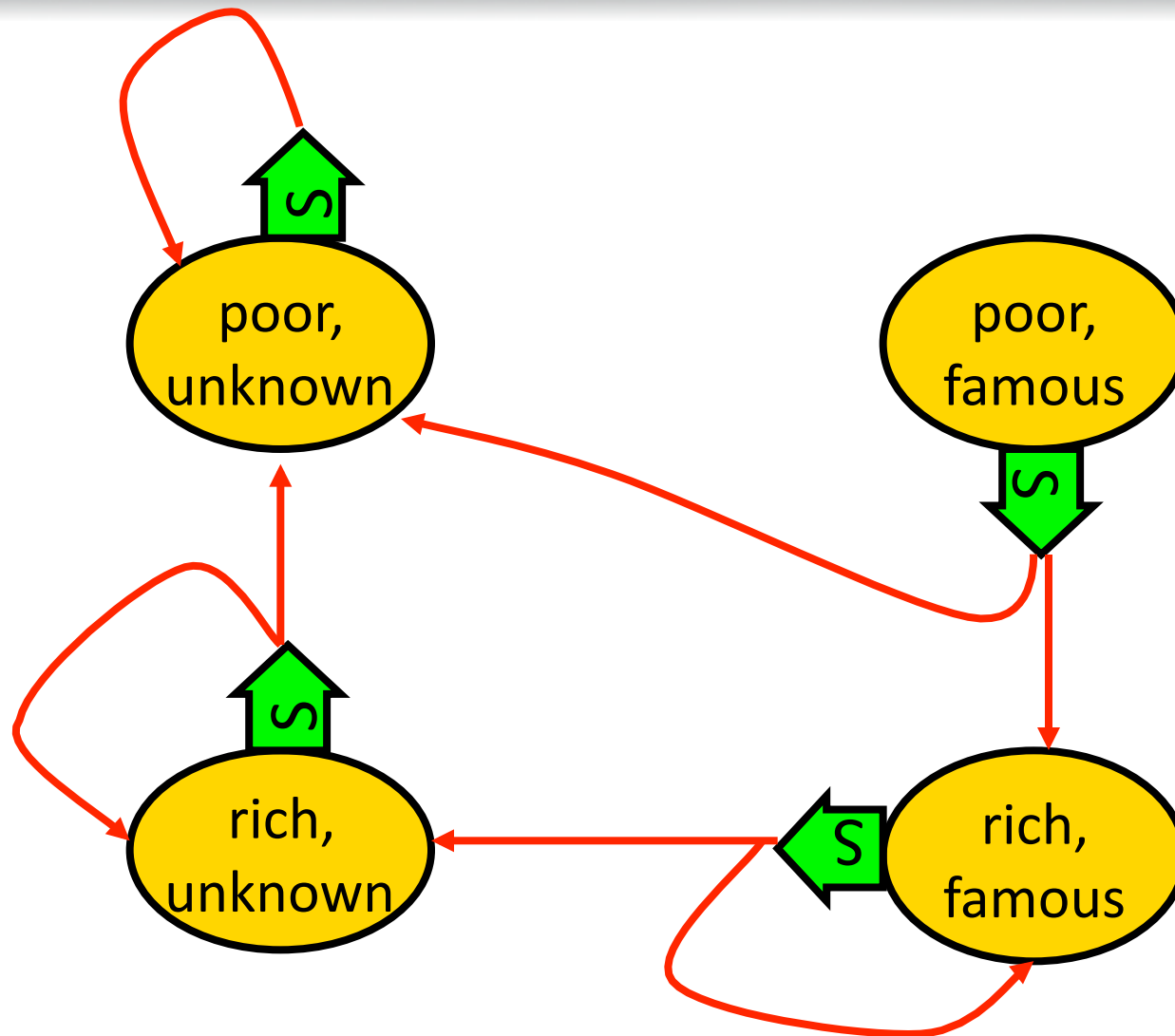- Fixed probability $(1-\gamma)$ of "obliteration" (inflation, running out of battery, …)

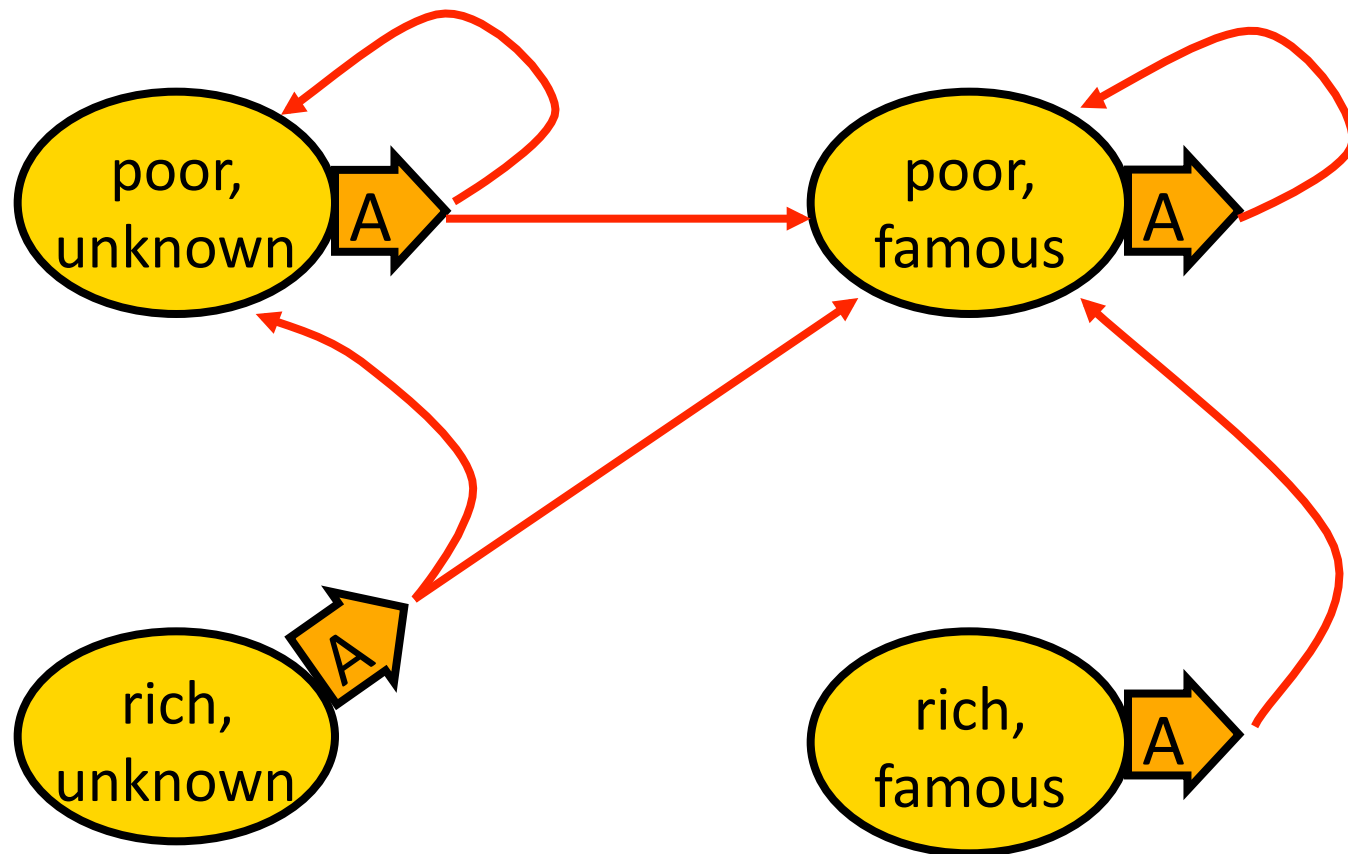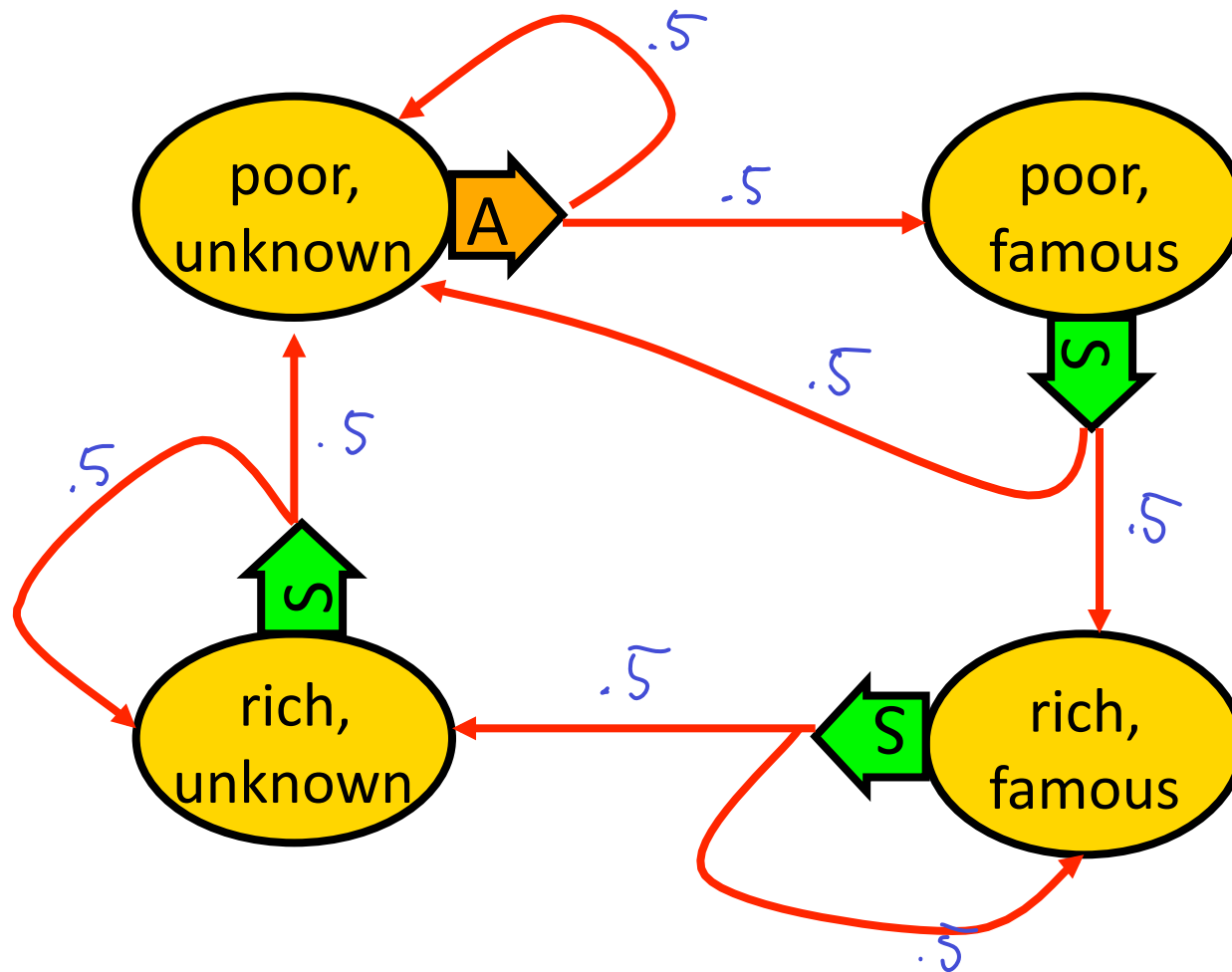**Policy**: Pick one *fixed* action for each state

24

# Planning in MDPs

- Deterministic policy $\quad\quad\quad \pi: X \to A$

- Induces a **Markov chain**: $\quad\quad X_1, X_2, \ldots, X_t, \ldots$
  with transition probabilities

$$P(X_{t+1}=x' \mid X_t=x) = P(x' \mid x, \pi(x))$$



- Expected value $J(\pi) = E[\quad r(X_1, \pi(X_1))$
  $\quad\quad\quad\quad\quad\quad\quad\quad + \gamma\, r(X_2, \pi(X_2))$
  $\quad\quad\quad\quad\quad\quad\quad\quad + \gamma^2\, r(X_3, \pi(X_3))$
  $\quad\quad\quad\quad\quad\quad\quad\quad + \ldots \quad\quad\quad\quad ]$

# Computing the value of a policy

- For fixed policy $\pi$ and each state x, define **value function**

$$V^\pi(x) = J(\pi \mid \text{start in state } x) = r(x,\pi(x)) + E[\sum_t \gamma^t\, r(X_t,\pi(X_t))]$$

Recursion: $V^\pi(x) = r(x,\pi(x)) + \gamma\, E\left[\sum_t \gamma^{t-1} r(X_t,\pi(X_t))\right]$

$$= r(x,\pi(x)) + \gamma \sum_{x'} P(x' \mid x,\pi(x))\, V^\pi(x')$$

and $J(\pi) = V^\pi(X_0)$
     $\uparrow$ start state

$[V^\pi(1) \dots V^\pi(n)]^T$    $[r(1,\pi(1)), \dots r(n,\pi(n))]^T$

In matrix notation: $V^\pi = r + \gamma\, T V^\pi$

$$\begin{pmatrix} P(1 \mid 1,\pi(1)) \dots P(n \mid 1,\pi(1)) \\ \vdots \\ P(1 \mid n,\pi(n)) \dots P(n \mid n,\pi(n)) \end{pmatrix}$$

$$\Rightarrow V^\pi = (I - \gamma T)^{-1} r$$

➔ **Can compute $V^\pi$ analytically, by matrix inversion!** ☺

**How can we find the optimal policy?**

# A simple algorithm

- For every policy $\pi$ compute $J(\pi)$
- Pick $\pi^* = \text{argmax } J(\pi)$

**Is this a good idea??**

No!!

# policies $= |A|^{|S|}$

Sps we start in state x

$$Q(x, a) = r(x, a) + \sum_{x'} P(x'|x, a) \cdot V(x')$$

$$\Rightarrow a^* \in \arg\max_a Q(x, a)$$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 3 | 0.812 | 0.868 | 0.918 | +1 |
| 2 | 0.762 | | 0.660 | −1 |
| 1 | 0.705 | 0.655 | 0.611 | 0.388 |

0.8

0.1          0.1

# Value functions and policies

Every value function induces a policy

**Value function $V^\pi$**

$V^\pi(x) = r(x,\pi(x)) +$
$\gamma\sum_{x'}P(x'|x,\pi(x)) V^\pi(x')$

**Greedy policy w.r.t. V**

$\pi_V(x) = \text{argmax}_a r(x,a)+$
$\gamma \sum_{x'} P(x' | x,a) V(x)$

Every policy induces a value function

**Thm**: Policy optimal $\Leftrightarrow$ greedy w.r.t. its induced value function!

# Policy iteration

- Start with a random policy $\pi$
- Until converged do:

  Compute value function $V_\pi(x)$

  Compute greedy policy $\pi_G$ w.r.t. $V_\pi$

  Set $\pi \leftarrow \pi_G$

- Guaranteed to
  - Monotonically improve
  - Converge to an optimal policy $\pi^*$

$$\forall t, x : V^{\pi_{t+1}}(x) \geq V^{\pi_t}(x)$$

- Often performs really well!
- Not known whether it's polynomial in |X| and |A|!

# Alternative approach

- For the optimal policy $\pi^*$ it holds **(Bellman equation)**

$$V^*(x) = \max_a r(x,a) + \gamma \sum_{x'} P(x' \mid x, a) V^*(x)$$

- Compute $V^*$ using dynamic programming:

$V_t(x) \quad = \qquad$ Max. expected reward when starting in state $x$ and world ends in $t$ time steps

$V_0(x) \quad = \quad \max\limits_a r(x,a)$

$V_1(x) \quad = \quad \max\limits_a r(x,a) + \gamma \cdot \sum_{x'} P(x' \mid x, a) V_0(x)$

$V_{t+1}(x) = \qquad \underline{\qquad} \quad '' \quad \underline{\qquad\qquad} \quad V_t(x)$

# Value iteration

- Initialize $V_0(x) = \max_a r(x,a)$
- For $t = 1$ to $1$

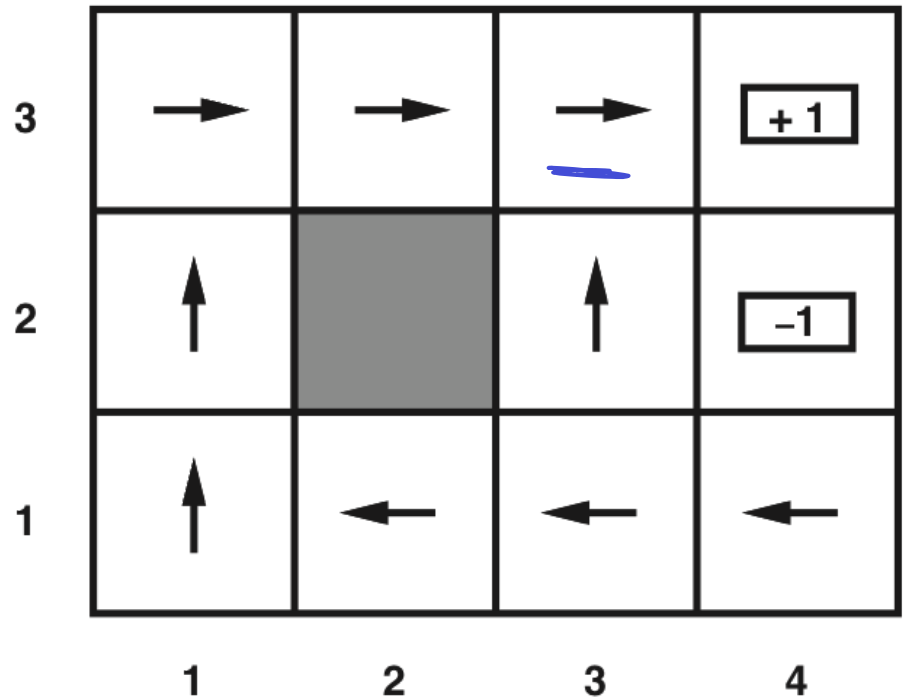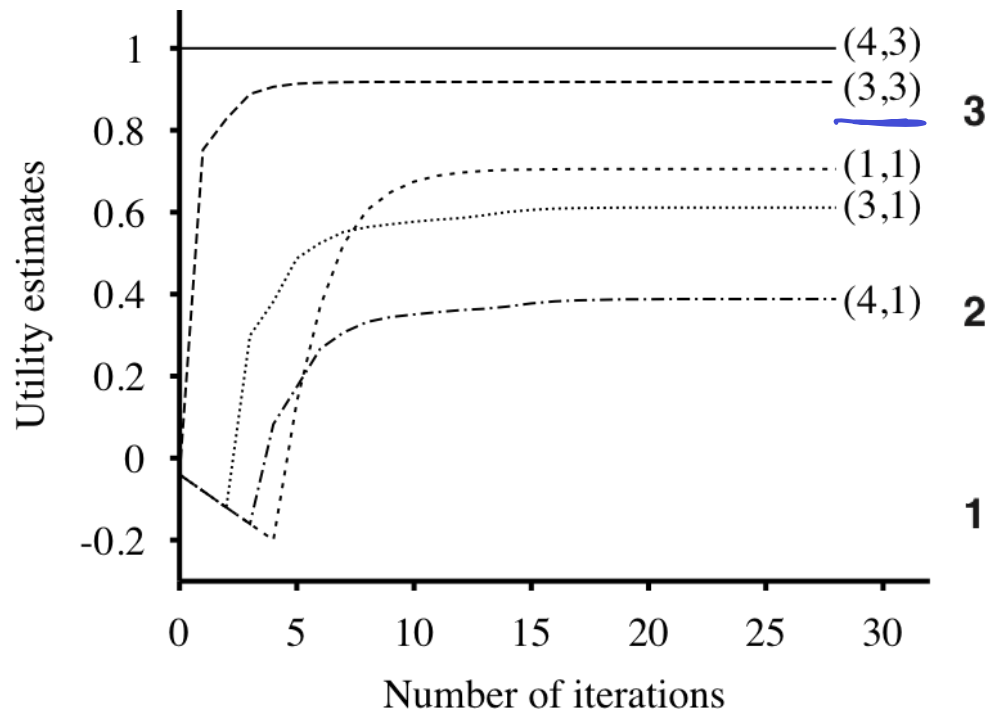  For each x, a, let $\quad Q_t(x,a) = r(x,a) + \gamma \sum_{x'} P(x'|x,a) V_{t-1}(x')$

  For each x let $\quad V_t(x) = \max_a Q_t(x,a)$

  Break if $\quad \max_x \left| V_t(x) - V_{t-1}(x) \right| \leq \varepsilon$

- Then choose greedy policy w.r.t. $V_t$

- **Guaranteed to converge to $\varepsilon$-optimal policy!**

# Value iteration

# Recap: Ways for solving MDPs

- Policy iteration:
  - Start with random policy $\pi$
  - Compute exact value function $V^\pi$ (matrix inversion)
  - Select greedy policy w.r.t. $V^\pi$ and iterate

- Value iteration
  - Solve Bellman equation using dynamic programming
    $V_t(x) = \max_a r(x,a) + \gamma \sum_{x`} P(x' \mid x,a) V_{t-1}(x)$

- Linear programming

# Applications of MDPs

- Robot path planning (noisy actions)

- Elevator scheduling

- Manufactoring processes

- Network switching and routing

- AI in computer games

- …