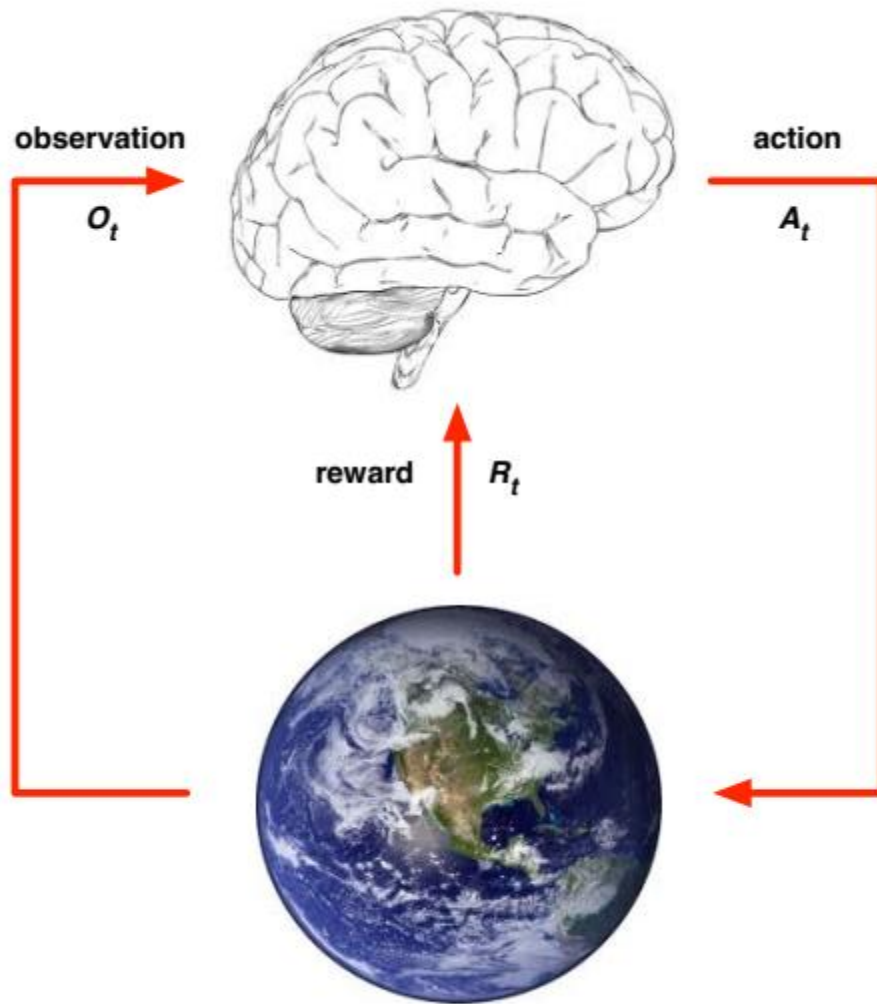


Lecture 1.1

# Reinforcement Learning: Markov Decision Processes

Alexey Gruzdev  
*alexey.s.gruzdev@gmail.com*  
HSE, Winter 2019

# Recap: RL Basics



An RL agent may include one or more of these components:

- **Policy:** agent's behavior function
- **Value function:** how good is each state and/or action
- **Model:** agent's representation of the environment

# Rewards hypothesis revisited

- A **reward**  $R_t$  is a scalar feedback signal
- Indicates how well agent is doing at step  $t$
- The agent's job is to maximize cumulative reward

Reinforcement Learning is based on the **reward hypothesis**.

The reward hypothesis:

*All* goals can be described by the maximization of expected cumulative reward .

- Online banners recommender system
- Personal mobile-phone assistants

# Markov Processes Family

- Markov Processes (Markov Chain)
- Markov Reward Processes
- Markov Decision Processes
- Extensions to MDPs:
  - Infinite & Continuous MDP
  - POMDP
  - Undiscounted MDP

# Introduction to MDPs

- *Markov decision processes* formally describe an environment for reinforcement learning
- Where the environment is *fully observable*
- i.e. The current *state* completely characterizes the process
- Almost all RL problems can be formalized as MDPs, e.g.
  - Optimal control primarily deals with continuous MDPs
  - Partially observable problems can be converted into MDPs
  - Bandits are MDPs with one state

# Markov Property

- Definition: a state  $S_t$  is **Markov** if and only if

$$P[S_{t+1} | S_t] = P[S_{t+1} | S_1, \dots, S_t]$$

- “The future is independent of the past given the present”
  - The state captures all relevant information from the history
  - Once the state is known, the history may be thrown away
  - i.e. The state is a sufficient statistic of the future

# State Transition Matrix

- For a Markov state  $s$  and successor state  $s'$ , the *state transition probability* is defined by

$$P_{ss'} = P [S_{t+1} = s' \mid S_t = s]$$

- State transition matrix  $\mathbf{P}_{ss'}$  defines transition probabilities from all states  $s$  to all successor states  $s'$

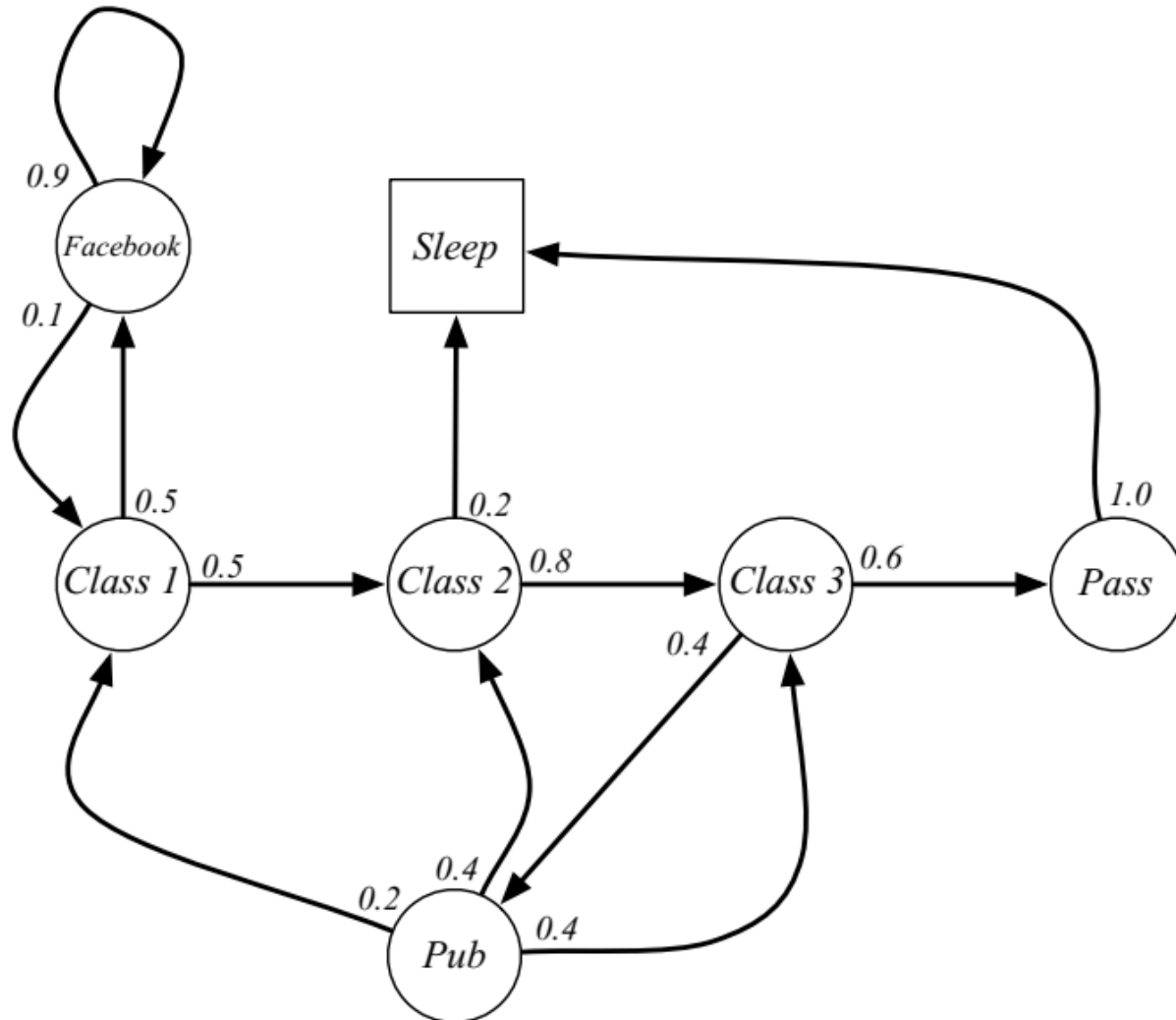
$$\mathbf{P}_{ss'} = \begin{pmatrix} P_{11} & \dots & P_{1n} \\ \vdots & \ddots & \vdots \\ P_{n1} & \dots & P_{nn} \end{pmatrix}$$

# Markov Process

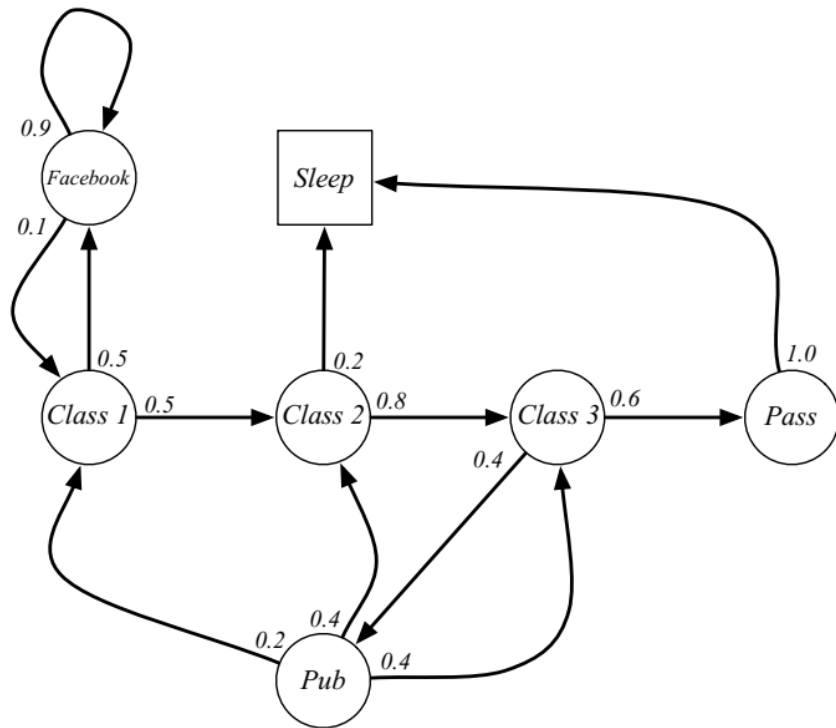
- Markov process is a memoryless random process, i.e. a sequence of random states  $S_1, \dots, S_t$  with the Markov property.
- *A Markov Process (or Markov Chain)* is a tuple  $(S, P)$ 
  - $S$  is a (finite) set of states
  - $P_{ss'}$  is a state transition probability matrix
  - $P_{ss'} = P[S_{t+1} = s' \mid S_t = s]$



# Markov Process: Student Example



# Markov Process: Episodes Sampling

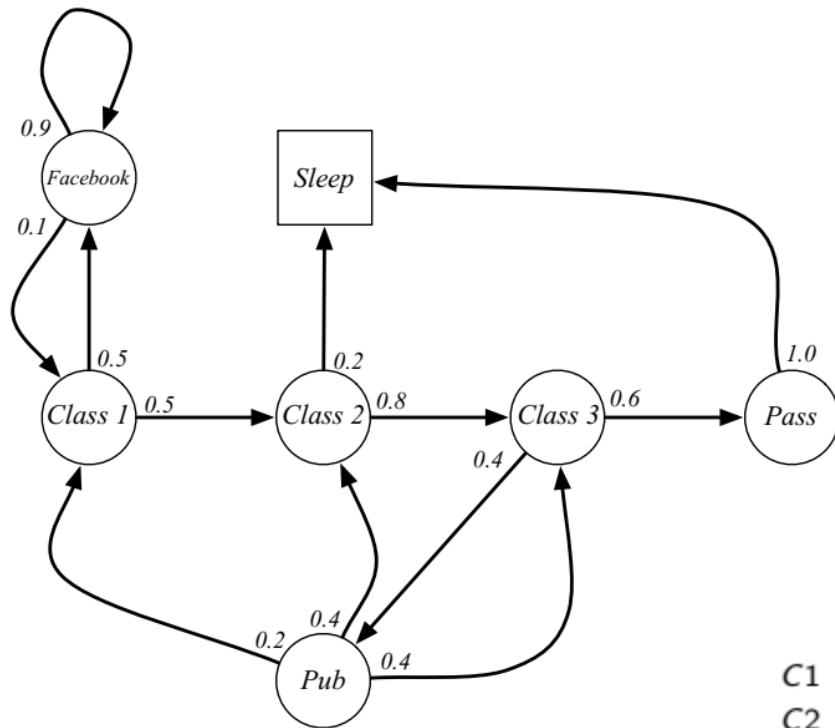


- Sample episodes for Student Markov Chain starting from  $S_1 = C1$

$S_1, \dots, S_t$

- C1 C2 C3 Pass Sleep
- C1 FB FB C1 C2 Sleep
- C1 C2 C3 Pub C2 C3 Pass Sleep
- C1 FB FB C1 C2 C3 Pub C1 FB  
FB FB C1 C2 C3 Pub C2 Sleep

# Markov Process: Transition Probabilities

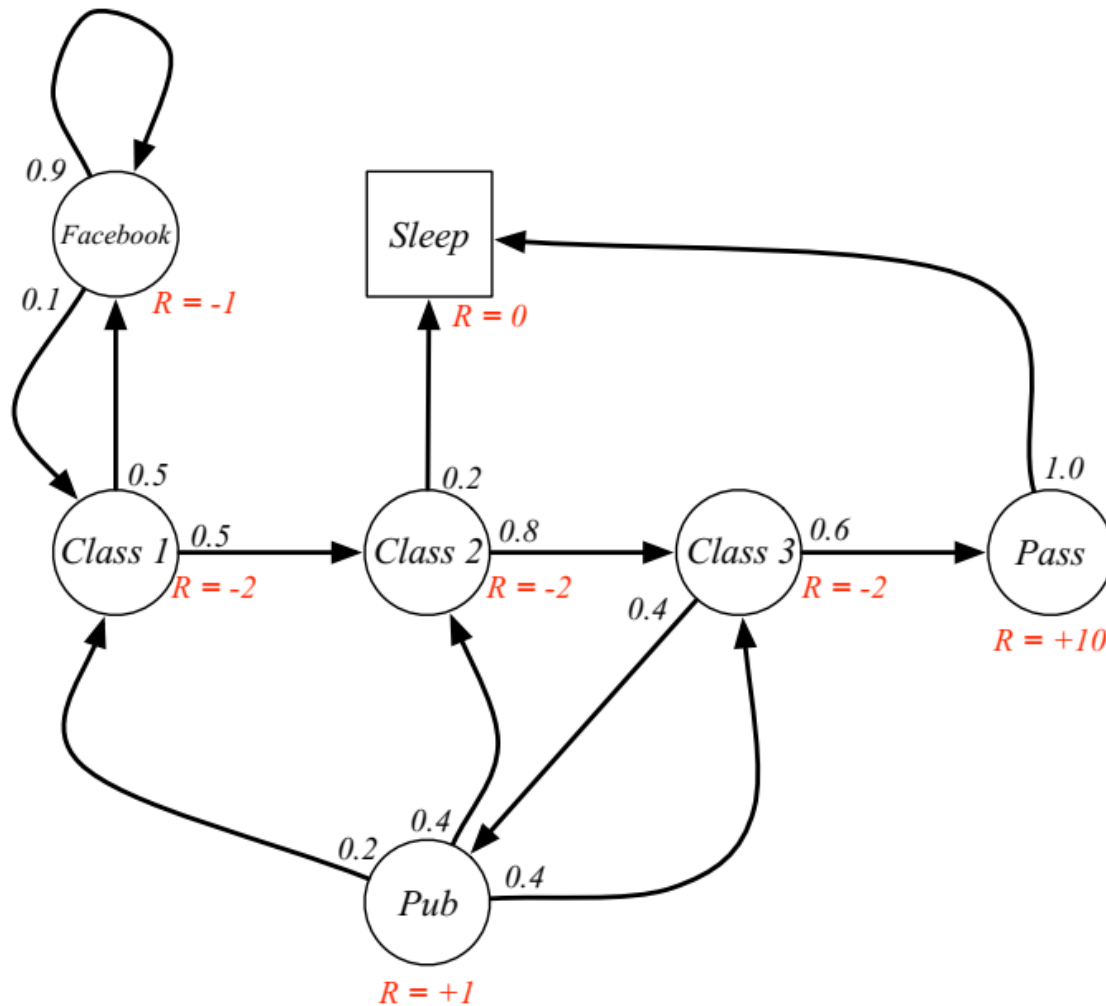


$$\mathcal{P} = \begin{matrix} & \begin{matrix} C1 & C2 & C3 & Pass & Pub & FB & Sleep \end{matrix} \\ \begin{matrix} C1 \\ C2 \\ C3 \\ Pass \\ Pub \\ FB \\ Sleep \end{matrix} & \begin{bmatrix} & & & & & 0.5 & \\ & 0.5 & & & & & 0.2 \\ & & 0.8 & & & & \\ & & & 0.6 & 0.4 & & 1.0 \\ 0.2 & 0.4 & 0.4 & & & & \\ 0.1 & & & & & 0.9 & \\ & & & & & & 1 \end{bmatrix} \end{matrix}$$

# Markov Reward Process

- A Markov reward process is a Markov chain with values.
- A *Markov Reward Process* is a tuple  $(S, P, R, \gamma)$ 
  - $S$  is a (finite) set of states
  - $P_{ss'}$  is a state transition probability matrix
  - $P_{ss'} = P[S_{t+1} = s' \mid S_t = s]$
  - $R$  is a reward function,  $R_s = E[R_{t+1} \mid S_t = s]$
  - $\gamma$  is a discount factor,  $\gamma \in [0, 1]$

# MRP: Student Example



# Return

- The *return*  $G_t$  is the total discounted reward from time-step  $t$ .

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

- The *discount*  $\gamma \in [0, 1]$  is the present value of future rewards
- The value of receiving reward  $R$  after  $k + 1$  time-steps is  $\gamma^k R$ .
- This values immediate reward above delayed reward.
  - $\gamma$  close to 0 leads to "myopic" evaluation
  - $\gamma$  close to 1 leads to "far-sighted" evaluation

# Discount intuition

Most Markov reward and decision processes are discounted. Why?

- Mathematically convenient to discount rewards
- Avoids infinite returns in cyclic Markov processes
- Uncertainty about the future may not be fully represented
- If the reward is financial, immediate rewards may earn more interest than delayed rewards
- Animal/human behavior shows preference for immediate reward
- It is sometimes possible to use *undiscounted* Markov reward processes (i.e.  $\gamma = 1$ ), e.g. if all sequences terminate

# Value Function

- The value function  $v(s)$  gives the long-term value of state  $s$
- The *state value function*  $v(s)$  of an MRP is the expected return starting from state  $s$

$$v(s) = E[G_t \mid S_t = s]$$



# Student MRP returns

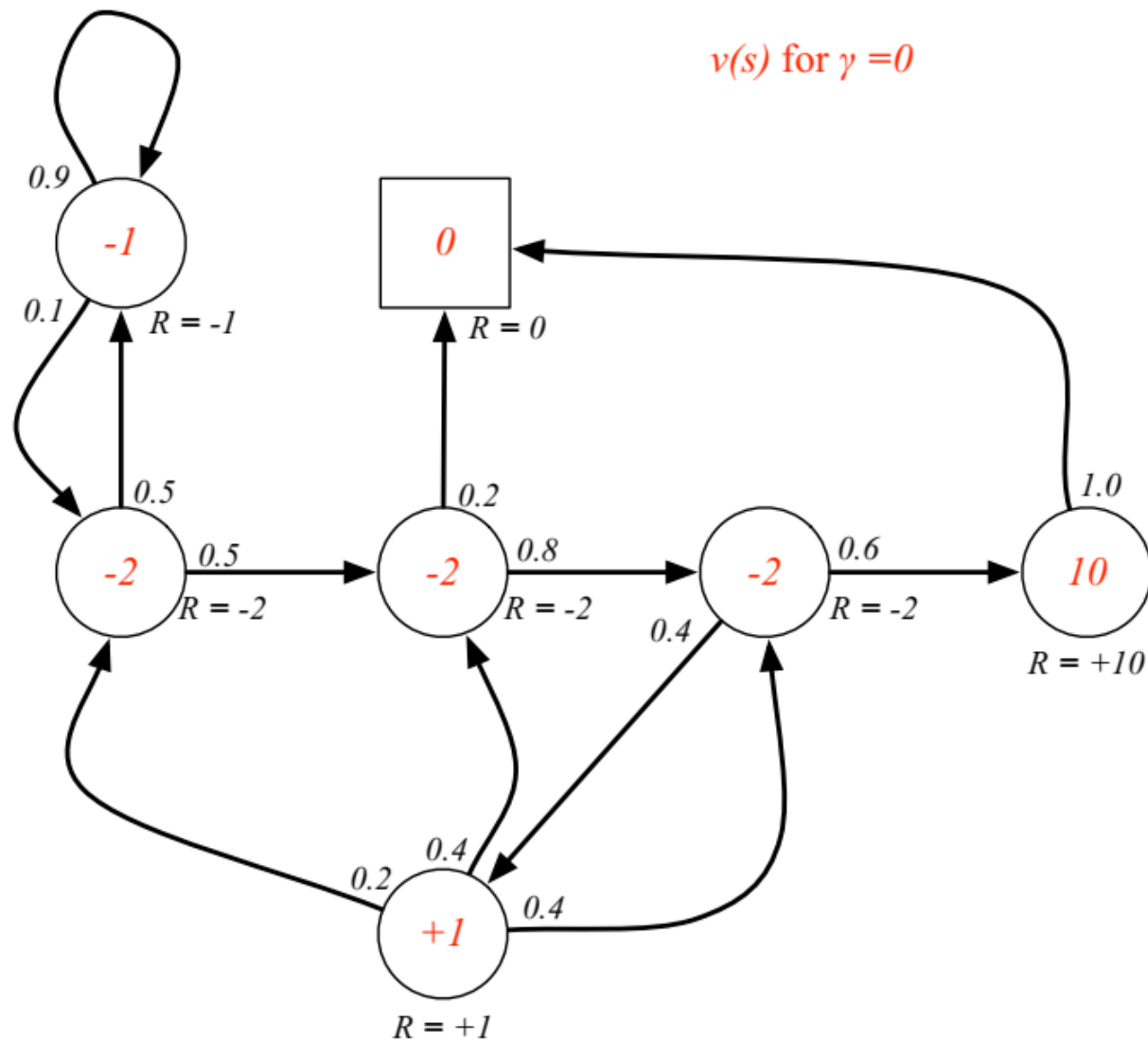
Sample returns for Student MRP with:

- $S_1 = C1$
- $\gamma = 0.5$

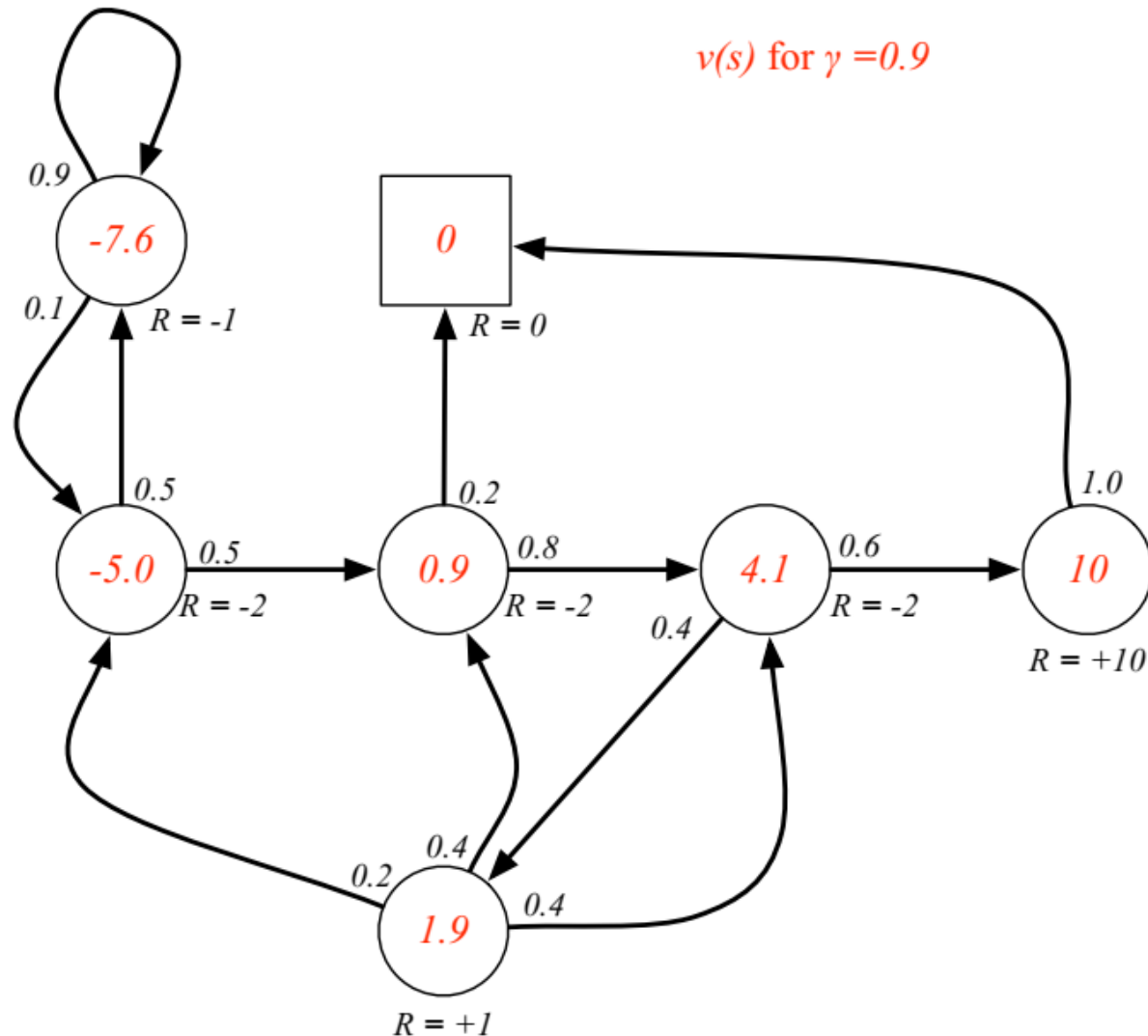
$$G_1 = R_2 + \gamma R_3 + \dots + \gamma^{T-2} R_T$$

C1 C2 C3 Pass Sleep	$v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 10 * \frac{1}{8}$	=	-2.25
C1 FB FB C1 C2 Sleep	$v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16}$	=	-3.125
C1 C2 C3 Pub C2 C3 Pass Sleep	$v_1 = -2 - 2 * \frac{1}{2} - 2 * \frac{1}{4} + 1 * \frac{1}{8} - 2 * \frac{1}{16} \dots$	=	-3.41
C1 FB FB C1 C2 C3 Pub C1 ...	$v_1 = -2 - 1 * \frac{1}{2} - 1 * \frac{1}{4} - 2 * \frac{1}{8} - 2 * \frac{1}{16} \dots$	=	-3.20
FB FB FB C1 C2 C3 Pub C2 Sleep			

# State-Value function for student MRP



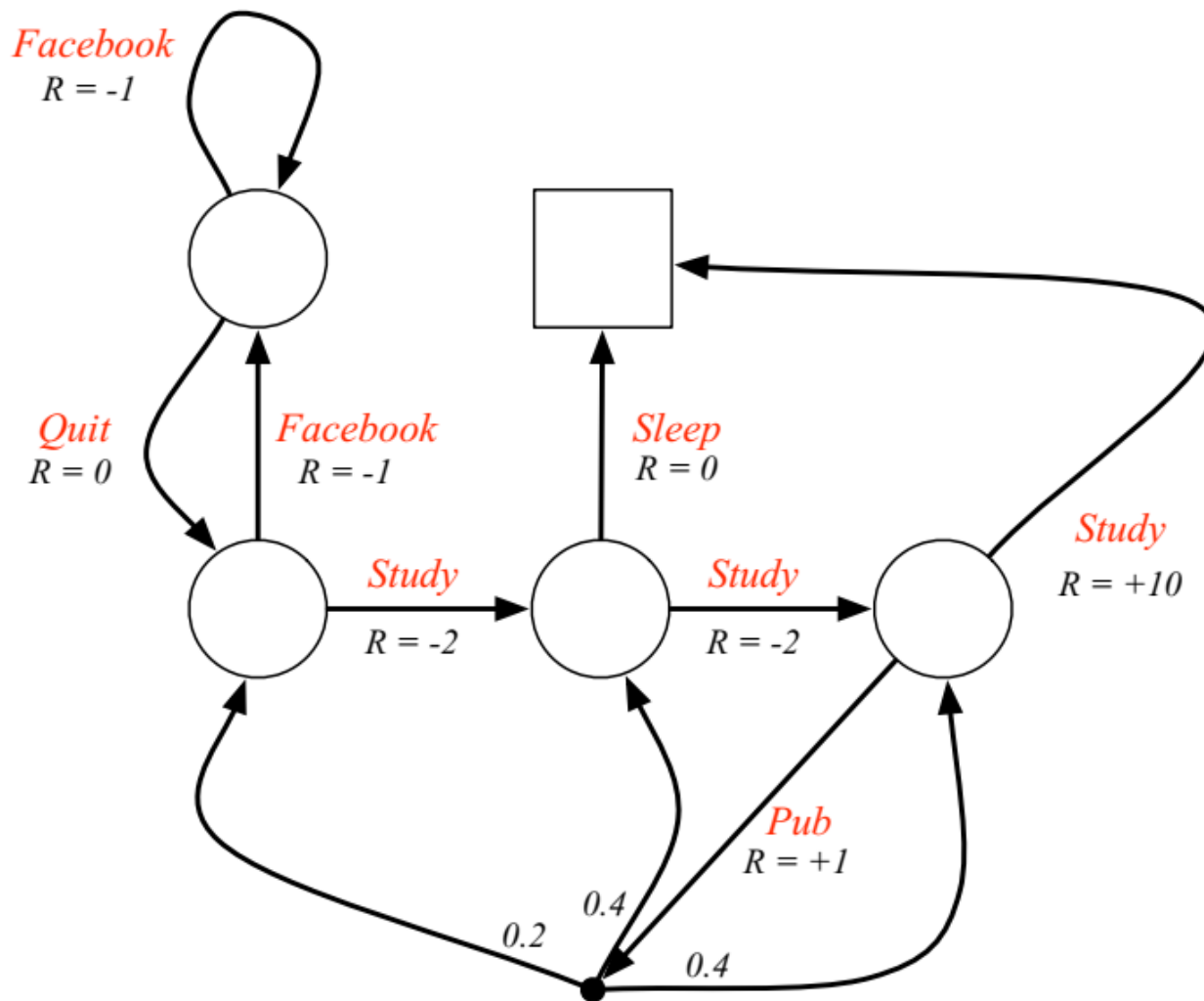
# State-Value function for student MRP



# Markov Decision Process

- A Markov decision process (MDP) is a Markov reward process with decisions. It is an *environment* in which all states are Markov.
- A Markov reward process is a Markov chain with values.
- A *Markov Decision Process* is a tuple  $(S, A, P, R, \gamma)$ 
  - $S$  is a (finite) set of states
  - $A$  is a finite set of actions
  - $P^a_{ss'}$  is a state transition probability matrix
  - $P^a_{ss'} = P[S_{t+1} = s' \mid S_t = s, A_t = a]$
  - $R$  is a reward function,  $R^a_s = E[R_{t+1} \mid S_t = s, A_t = a]$
  - $\gamma$  is a discount factor,  $\gamma \in [0, 1]$

# Markov Decision Process: Example



# MDP Policies

A *policy*  $\pi$  is a distribution over actions given states

$$\pi(a \mid s) = P[A_t = a \mid S_t = s]$$

- A policy fully defines the behavior of an agent
- MDP policies depend on the current state (not the history)
- i.e. Policies are *stationary* (time-independent),  
 $A_t \sim \pi(\cdot \mid S_t); \forall t > 0$

# MDP Policies

- Given an MDP  $M = (\mathcal{S}, \mathcal{A}, P, R, \gamma)$  and a policy  $\pi$
- The state sequence  $S_1, \dots, S_t$  is a Markov process  $(S, P^\pi)$
- The state and reward sequence  $S_1, R_1, S_2, \dots$  is a Markov reward process  $(S, P^\pi, R^\pi, \gamma)$
- where

$$P^\pi_{s,s'} = \sum_{a \in \mathcal{A}} \pi(a | s) P^a_{s,s'}$$

$$R^\pi_s = \sum_{a \in \mathcal{A}} \pi(a | s) R^a_s$$

# Updated Value functions

- The *state-value function*  $v_{\pi}(s)$  of an MDP is the expected return starting from state  $s$ , and then following policy  $\pi$

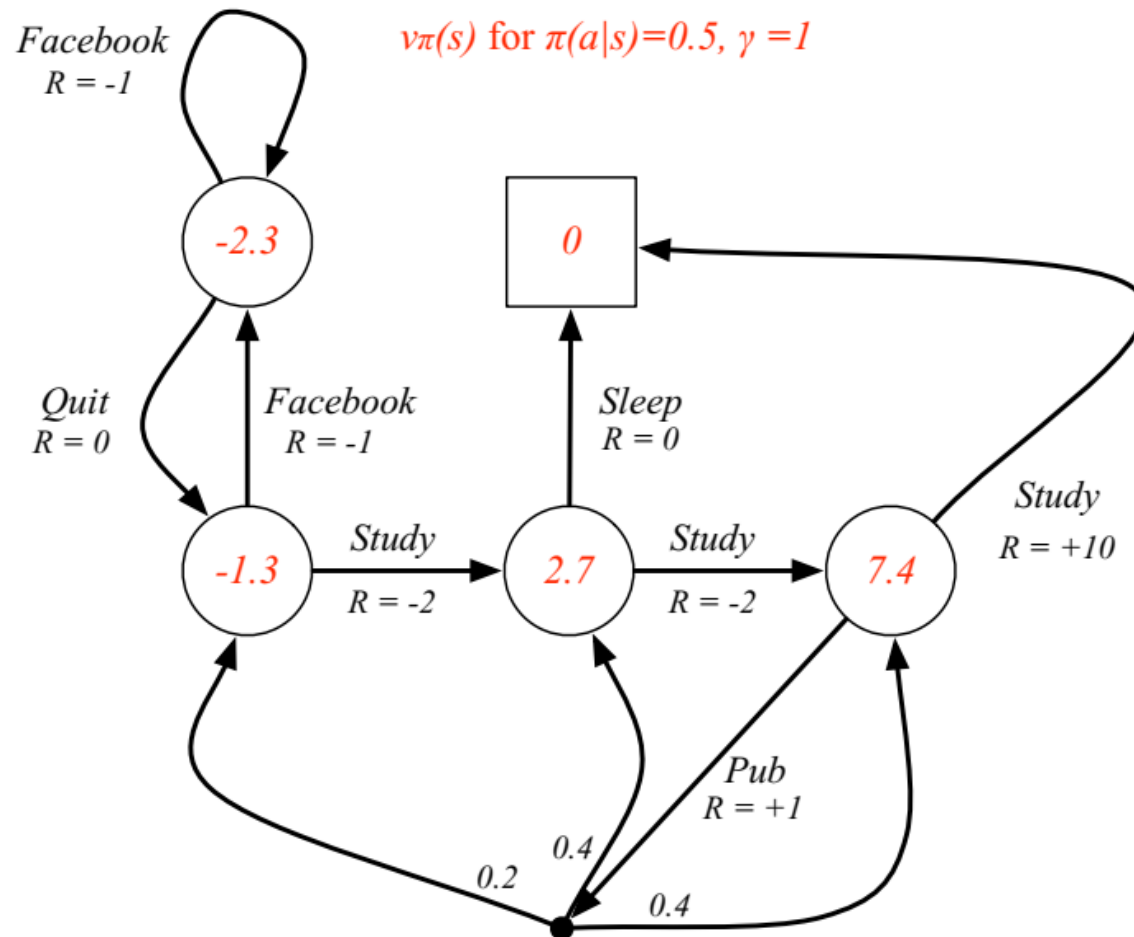
$$v_{\pi}(s) = E_{\pi} [G_t \mid S_t = s]$$

- The *action-value function*  $q_{\pi}(s, a)$  is the expected return starting from state  $s$ , taking action  $a$ , and then following policy  $\pi$

$$q_{\pi}(s, a) = E_{\pi} [G_t \mid S_t = s, A_t = a]$$



# Student MDP: Value function



# Optimal Value Function

- The *optimal state-value function*  $v_*(s)$  is the maximum value function over all policies

$$v_*(s) = \max_{\pi} (v_{\pi}(s))$$

- The *optimal action-value function*  $q_*(s; a)$  is the maximum action-value function over all policies

$$q_*(s, a) = \max_{\pi} (q_{\pi}(s, a))$$

- The optimal value function specifies the best possible performance in the MDP.
- An MDP is “solved” when we know the optimal value fn.

# Optimal Policy

- Define a partial ordering over policies:

$$\pi \geq \pi' \text{ if } v_{\pi}(s) \geq v_{\pi'}(s) \quad \forall s$$

**Theorem:** *For any Markov Decision Process*

- *There exists an optimal policy  $\pi_*$  that is better than or equal to all other policies:  $\pi_* \geq \pi \quad \forall \pi$*
- *All optimal policies achieve the optimal value function,*

$$v_{\pi_*}(s) = v_*(s)$$

- *All optimal policies achieve the optimal action-value function*

$$q_{\pi_*}(s, a) = q_*(s, a)$$

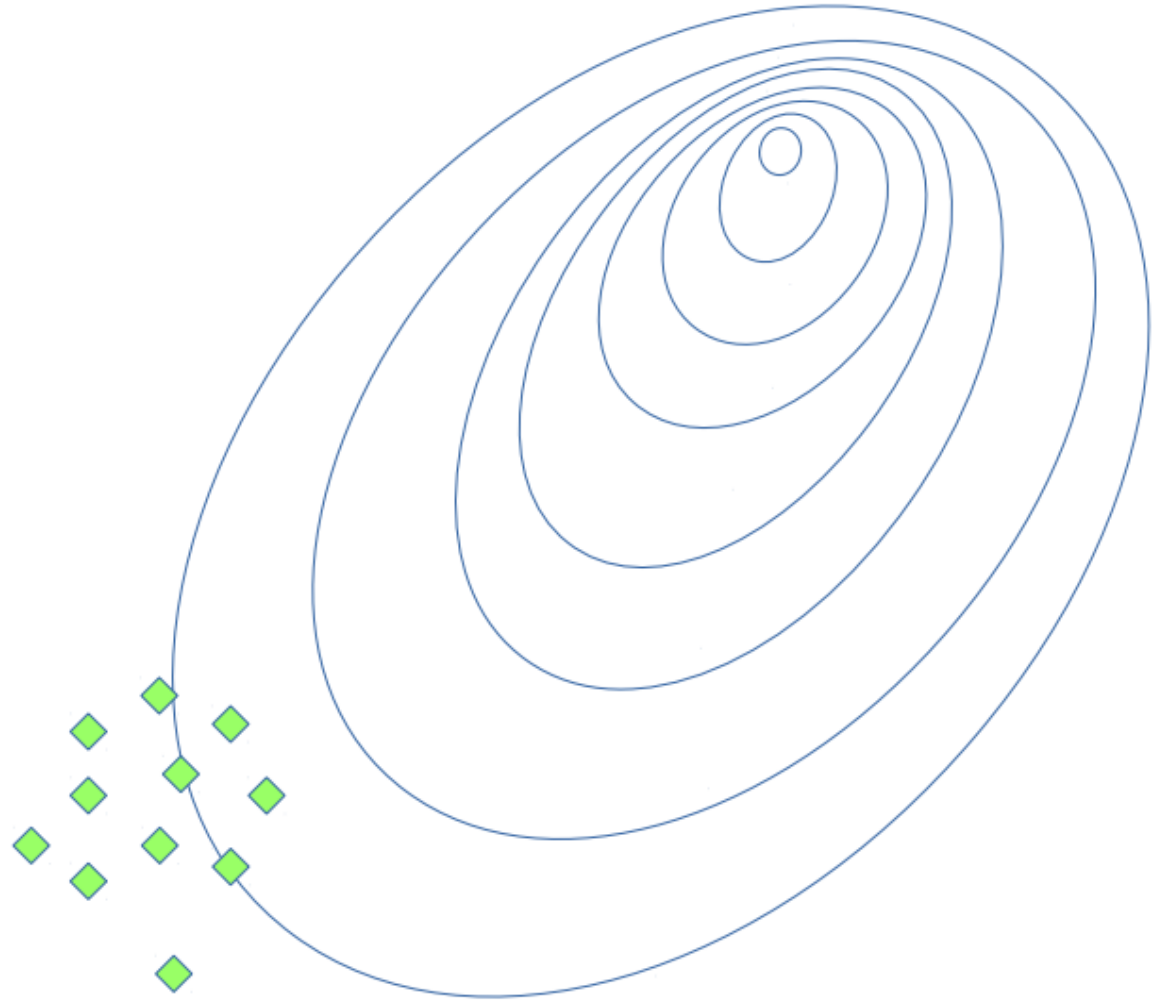
# No more theory – let's do real things!



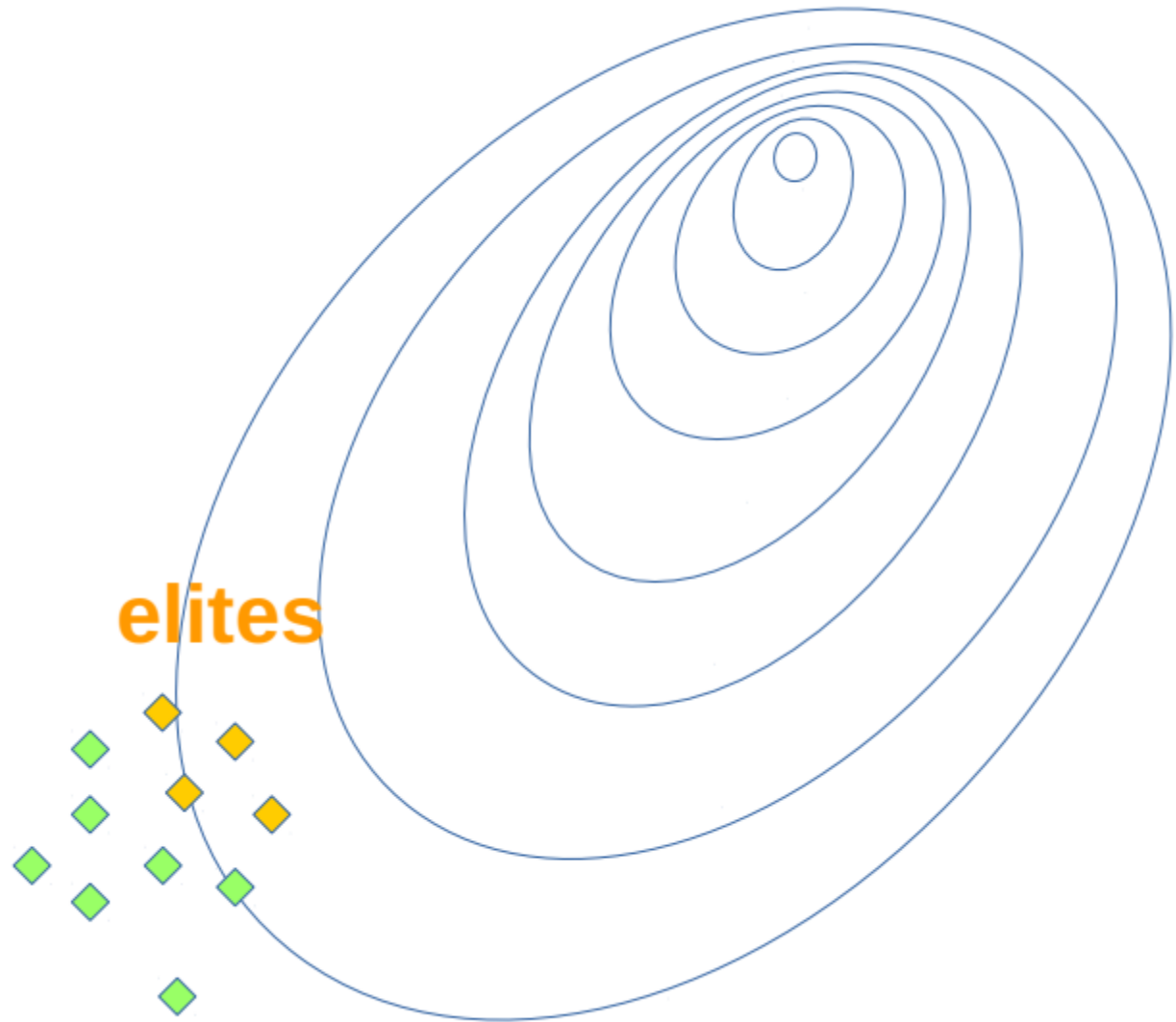
# Cross Entropy method

- Let's derive some heuristic:
  - Initialize policy (random!)
  - Repeat:
    - Sample N episodes
    - Pick best 30 % best episodes – a.k.a. “elite” episodes
    - Change policy to choose actions from elite episodes

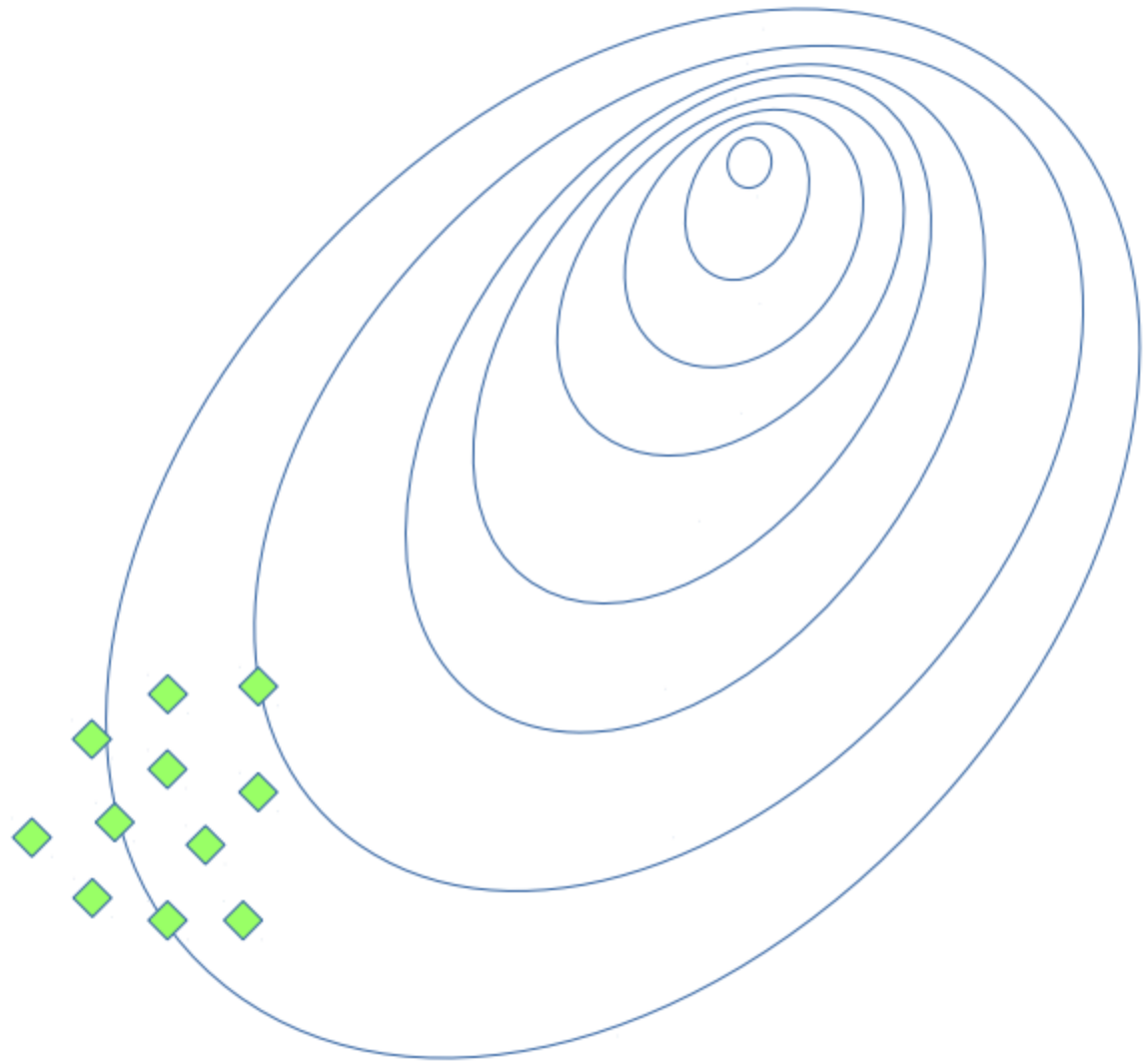
# Cross Entropy method: step-by-step



# Cross Entropy method: step-by-step

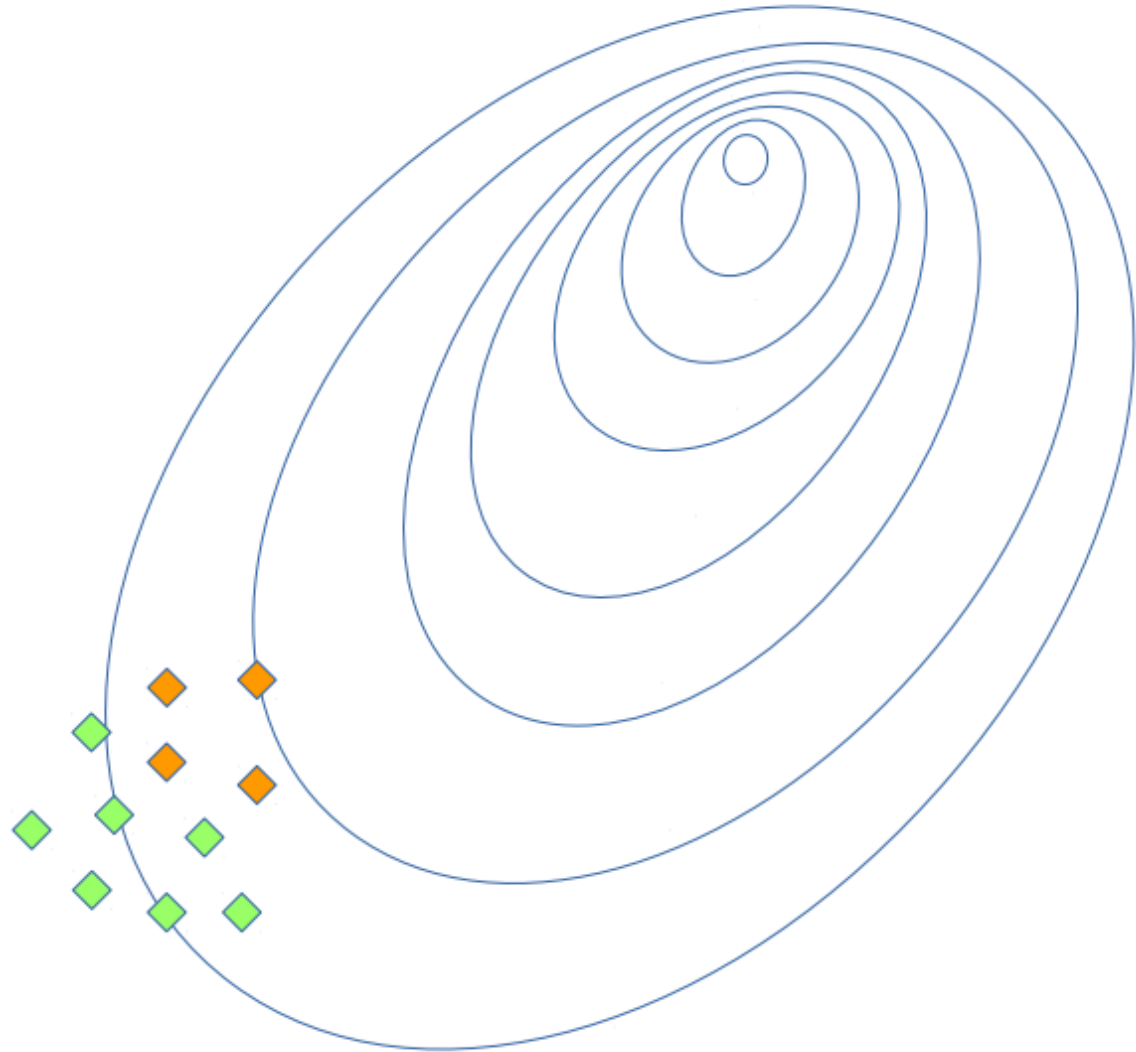


# Cross Entropy method: step-by-step

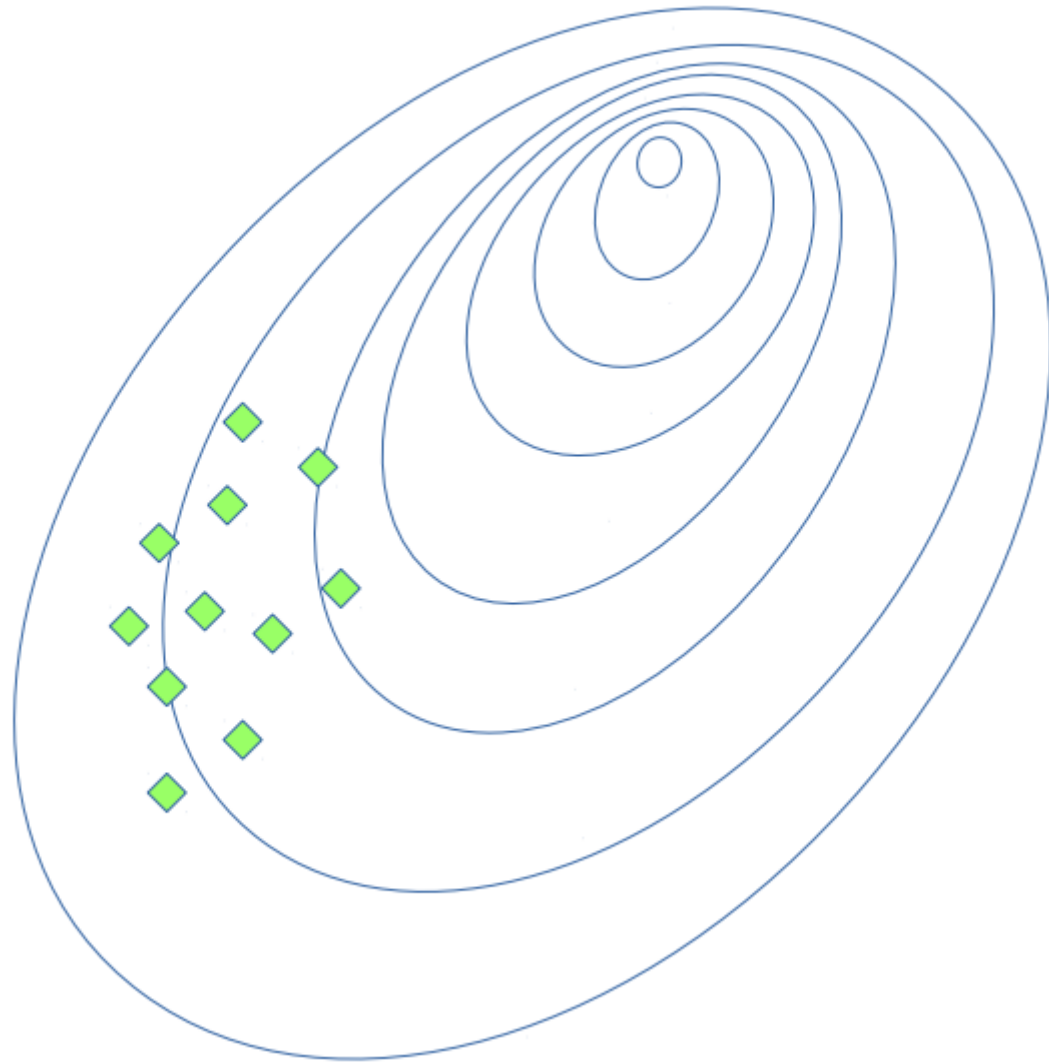




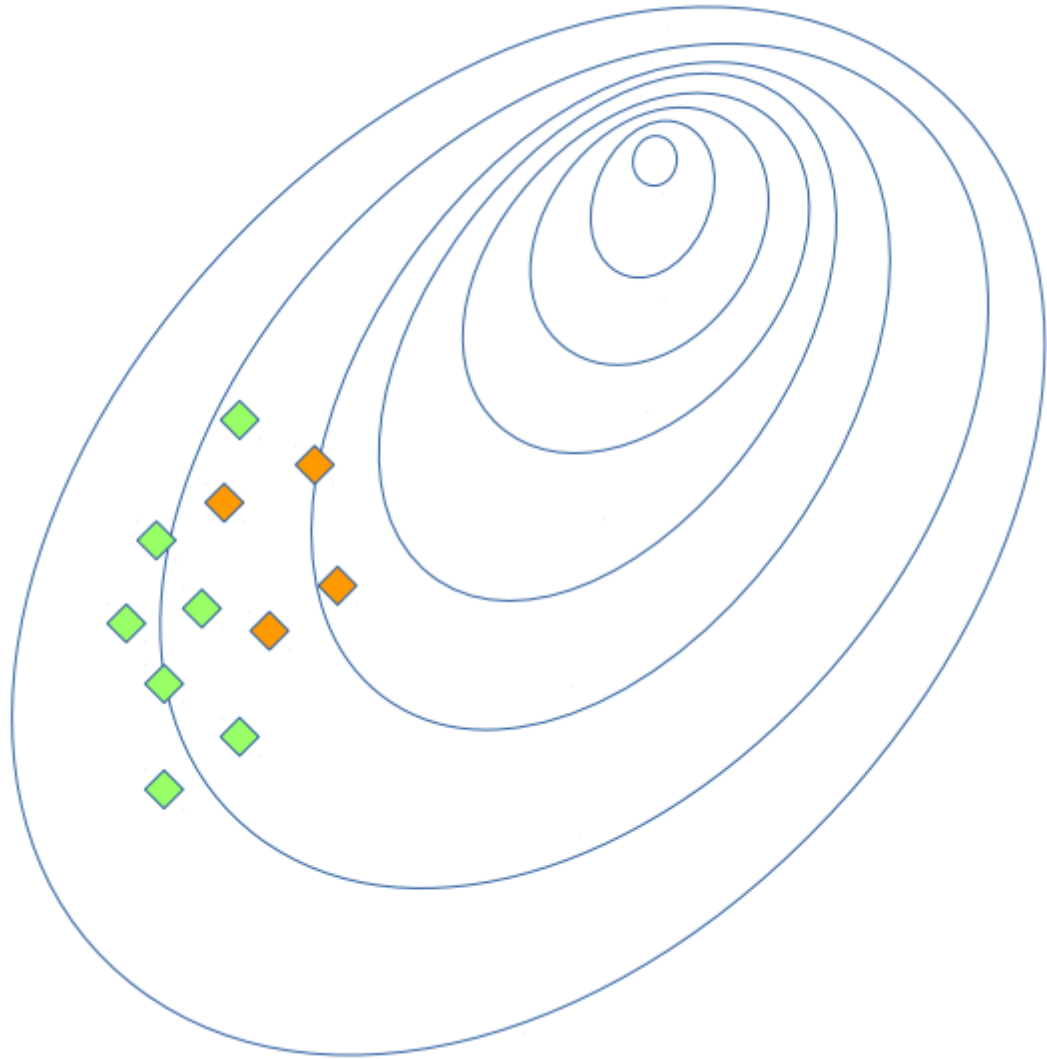
# Cross Entropy method: step-by-step



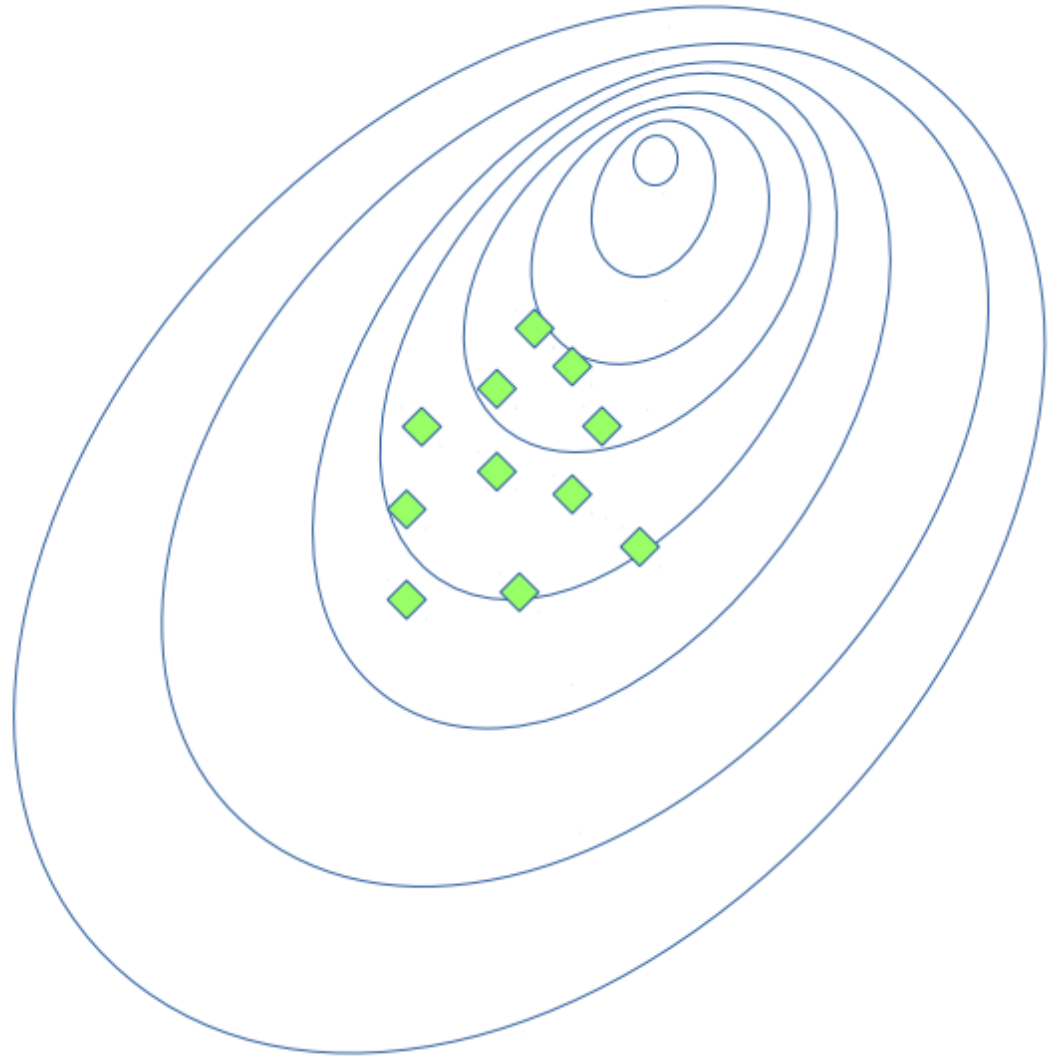
# Cross Entropy method: step-by-step



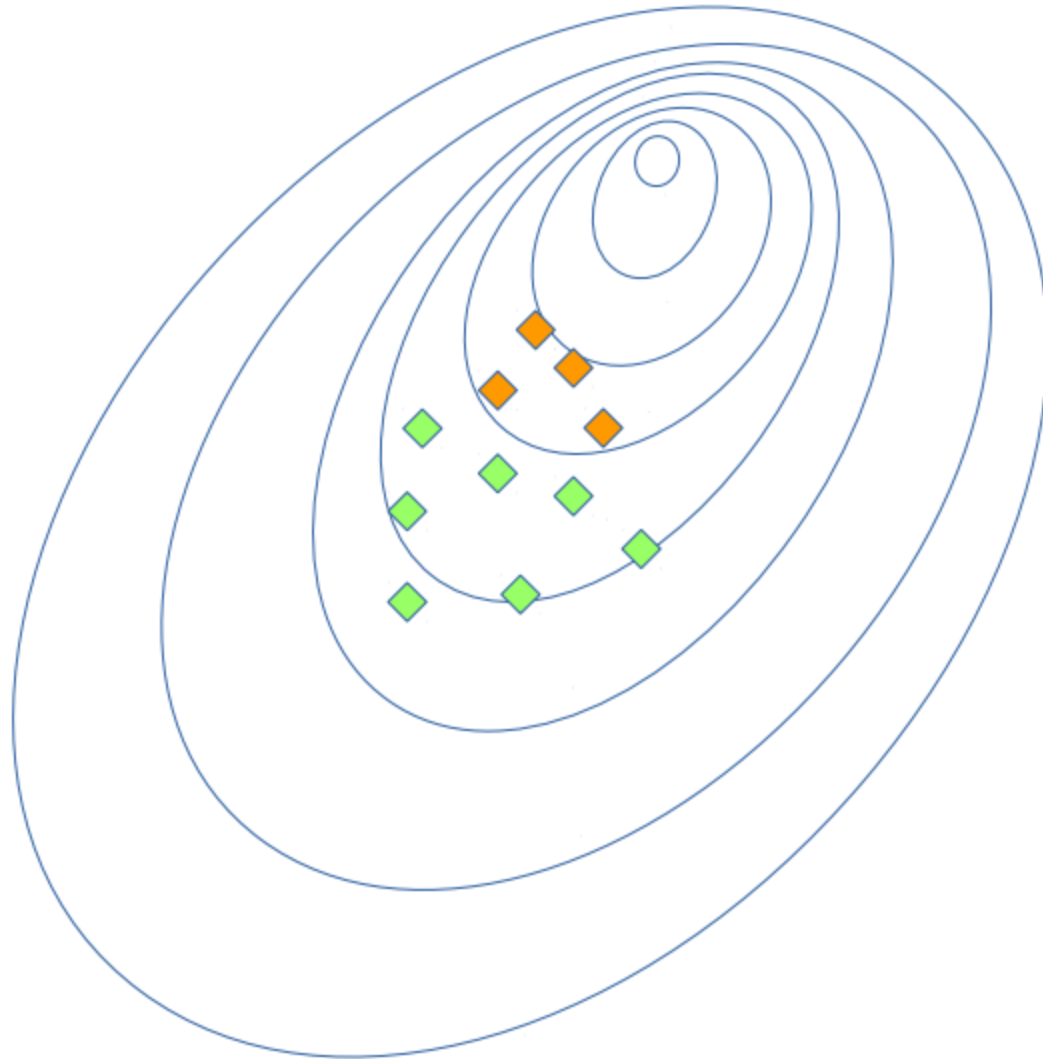
# Cross Entropy method: step-by-step



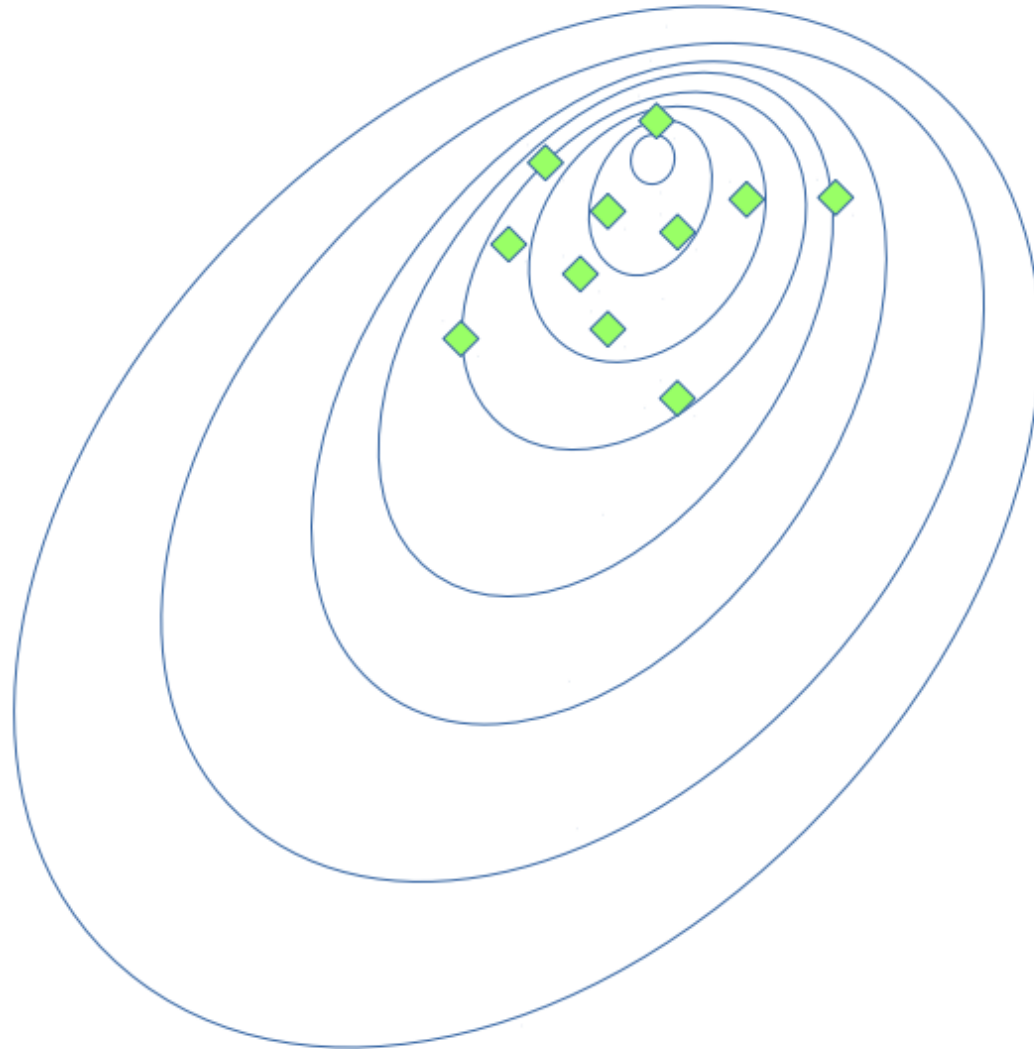
# Cross Entropy method: step-by-step



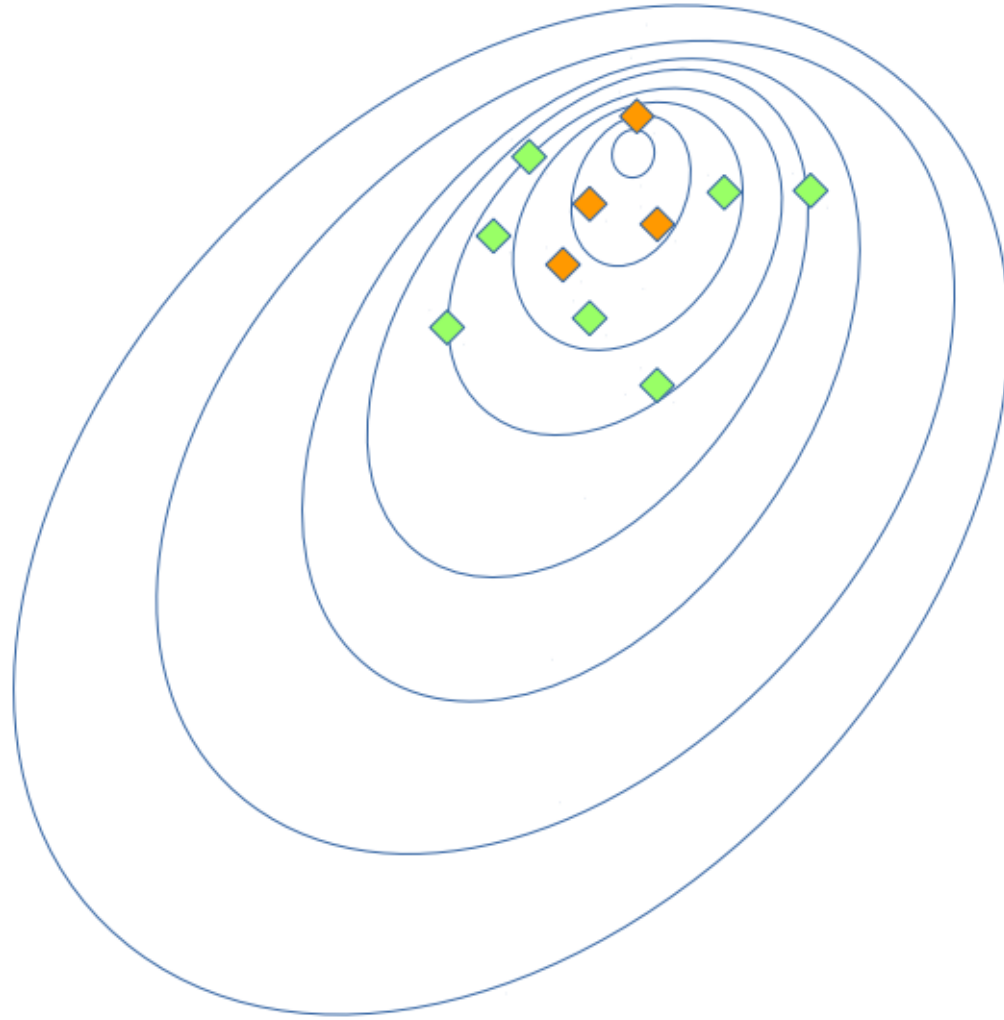
# Cross Entropy method: step-by-step



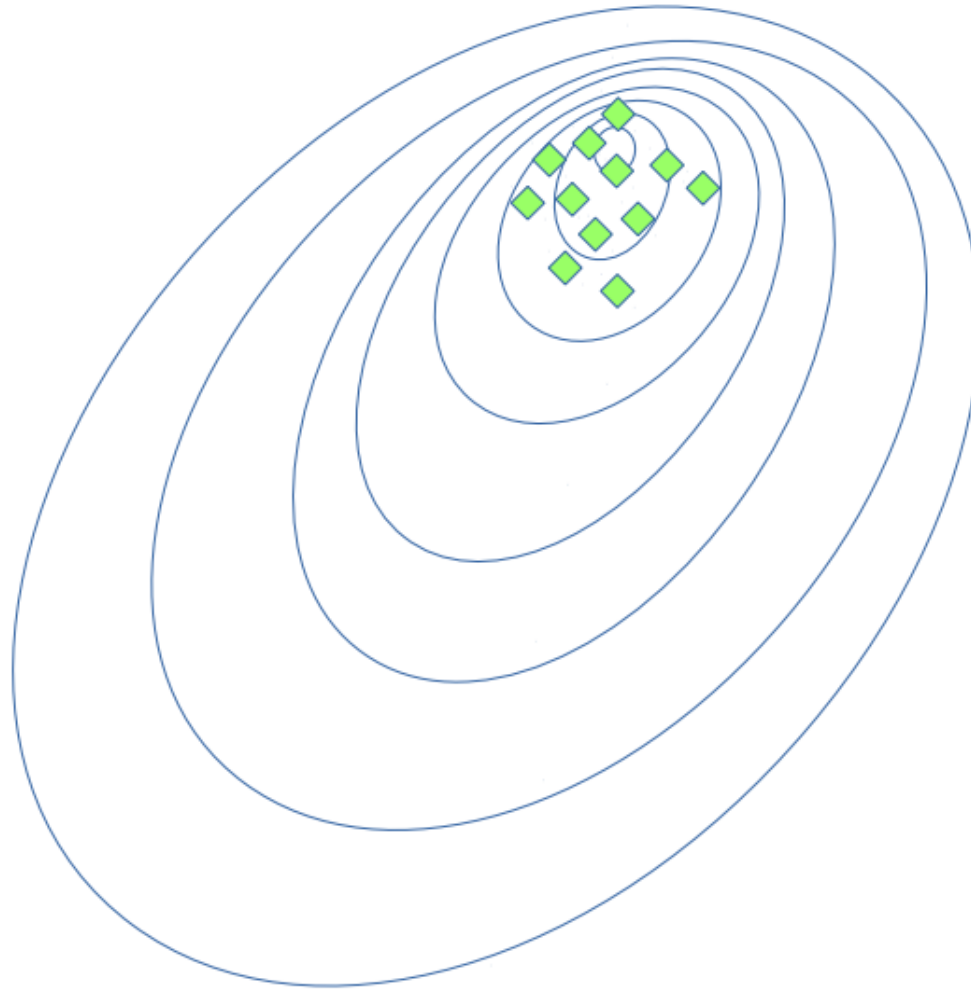
# Cross Entropy method: step-by-step



# Cross Entropy method: step-by-step



# Cross Entropy method: step-by-step





# Tabular Cross Entropy method

- Policy is matrix  $A$ 
  - $\pi(a \mid s) = P[A_t = a \mid S_t = s] = A_{s,a}$
- Sample  $N$  sessions with that policy
- Get  $M$  best sessions (elites)
- Elite =  $[(s_1, a_1), (s_2, a_2), \dots, (s_k, a_k)]$
- Update policy:

$$\pi(a \mid s) = \frac{\text{took } a \text{ at } s \text{ state}}{\text{was at } s \text{ state}} = \frac{\sum [s_t = s][a_t = a]}{\sum [s_t = s]}$$

# Cross Entropy problems

*But what if your environment has infinite/large state space ?*



# Approximated Cross Entropy method

- Policy is approximated
  - $\pi(a | s)$  predicted by Neural Network, Random Forest or any other ML algorithm.
- You can't set  $\pi(a | s)$  explicitly, it's not matrix anymore.
- Training data for our model:

$$\text{Elite} = [(s_1, a_1), (s_2, a_2), \dots, (s_k, a_k)]$$

# Questions?

