

# 10-701: Introduction to Machine Learning

## Lecture 16: Value and Policy Iteration

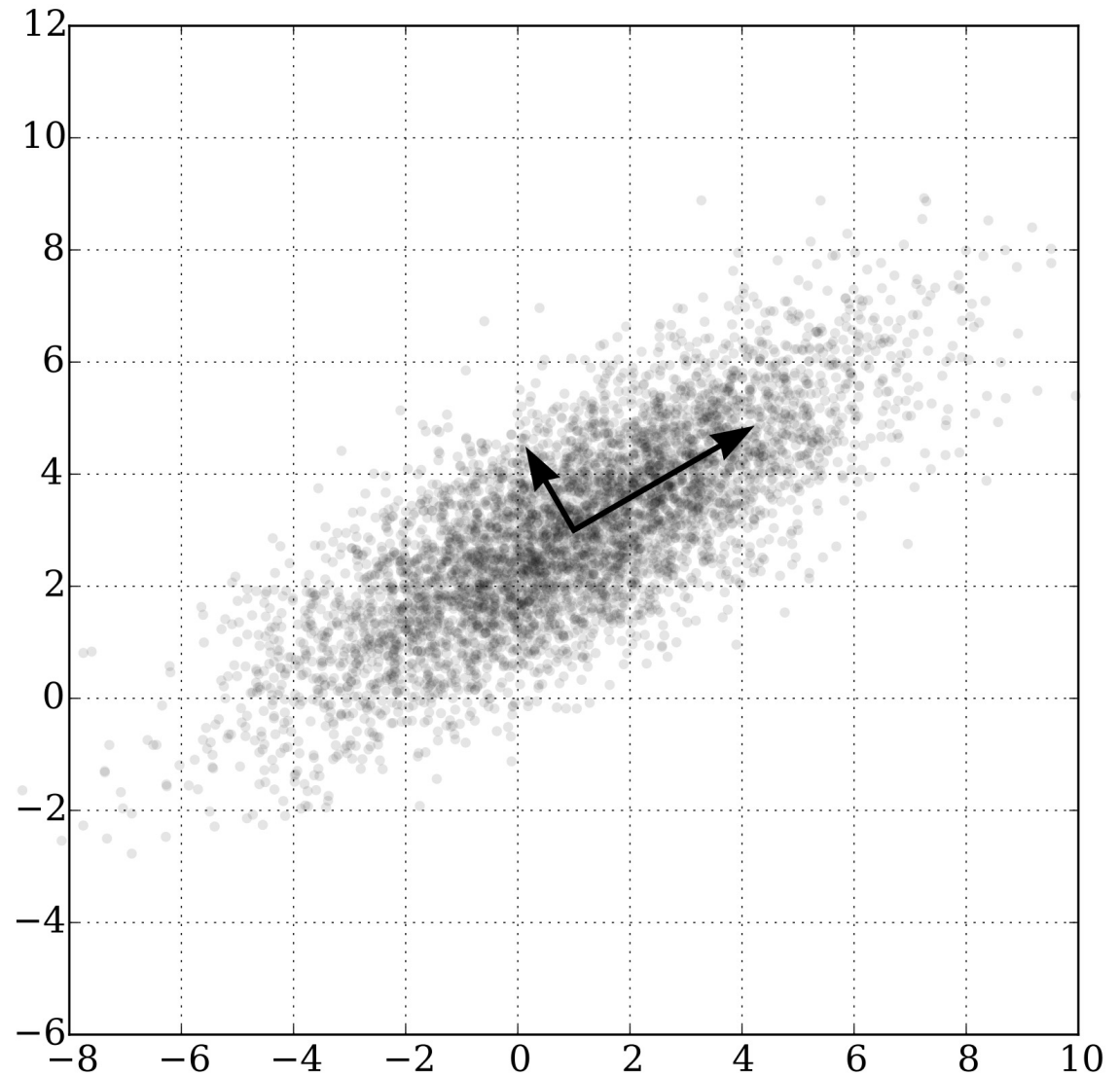
Henry Chai

3/18/24

# Front Matter

- Announcements
  - Midterm exam on 3/19 (tomorrow!) from **7 – 9 PM in DH A302**
  - Project proposals due on 3/22 (Friday) at 11:59 PM
    - You should submit proposals as a group, not individually: each group only needs to submit a single PDF
  - HW5 released 3/22 (Friday), due 4/1 at 11:59 PM
    - This is a *shorter, written-only HW*; you are expected to be working on your projects concurrently
- Recommended Readings
  - Mitchell, [Chapter 13](#)

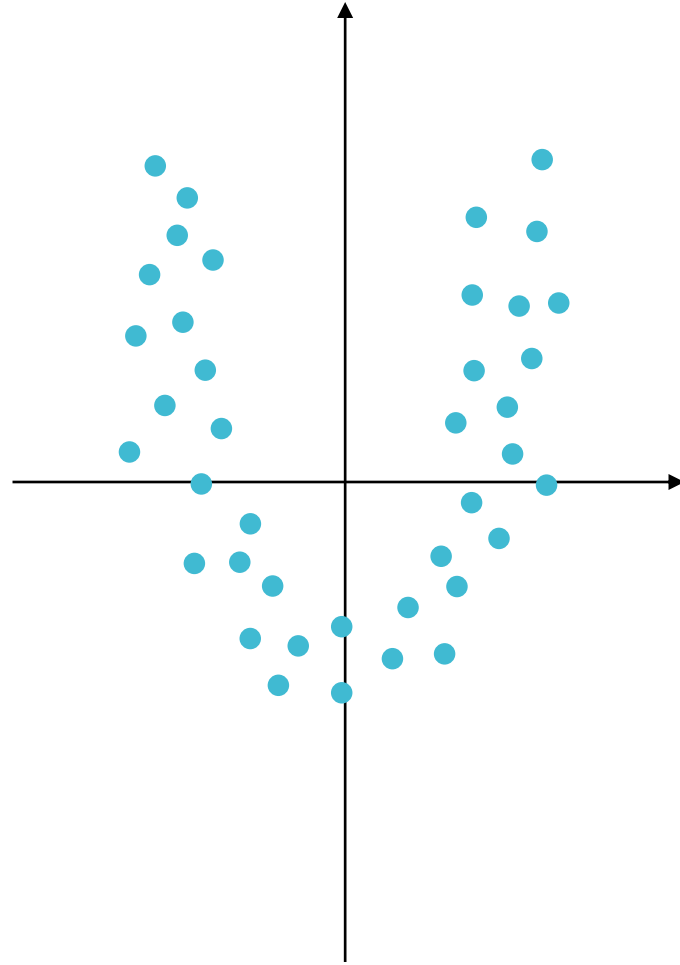
# Principal Components: Example



# PCA Algorithm

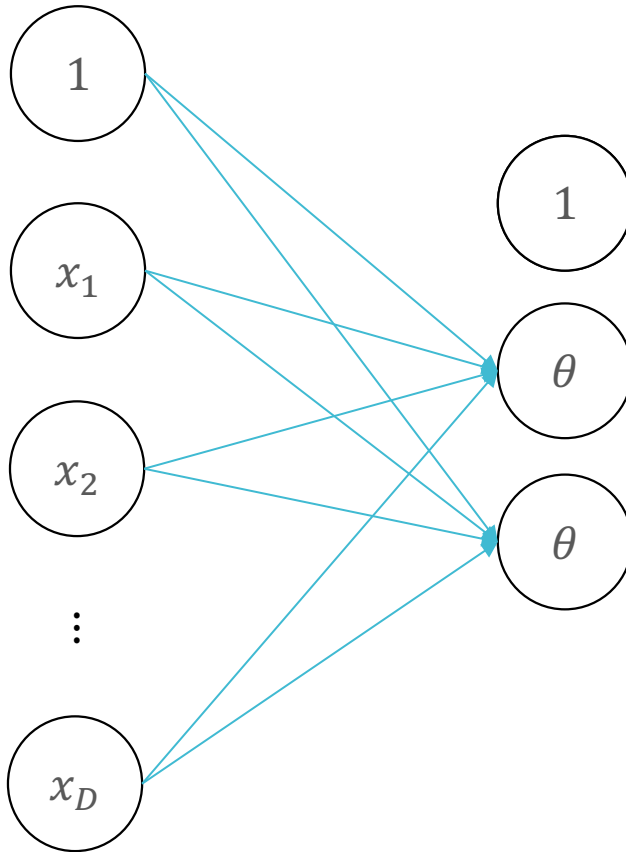
- Input:  $\mathcal{D} = \{(\mathbf{x}^{(n)})\}_{n=1}^N, \rho$ 
  1. Center the data
  2. Use SVD to compute the eigenvalues and eigenvectors of  $X^T X$
  3. Collect the top  $\rho$  eigenvectors (corresponding to the  $\rho$  largest eigenvalues),  $V_\rho \in \mathbb{R}^{D \times \rho}$
  4. Project the data into the space defined by  $V_\rho$ ,  $Z = X V_\rho$
- Output:  $Z$ , the transformed (potentially lower-dimensional) data

# Shortcomings of PCA



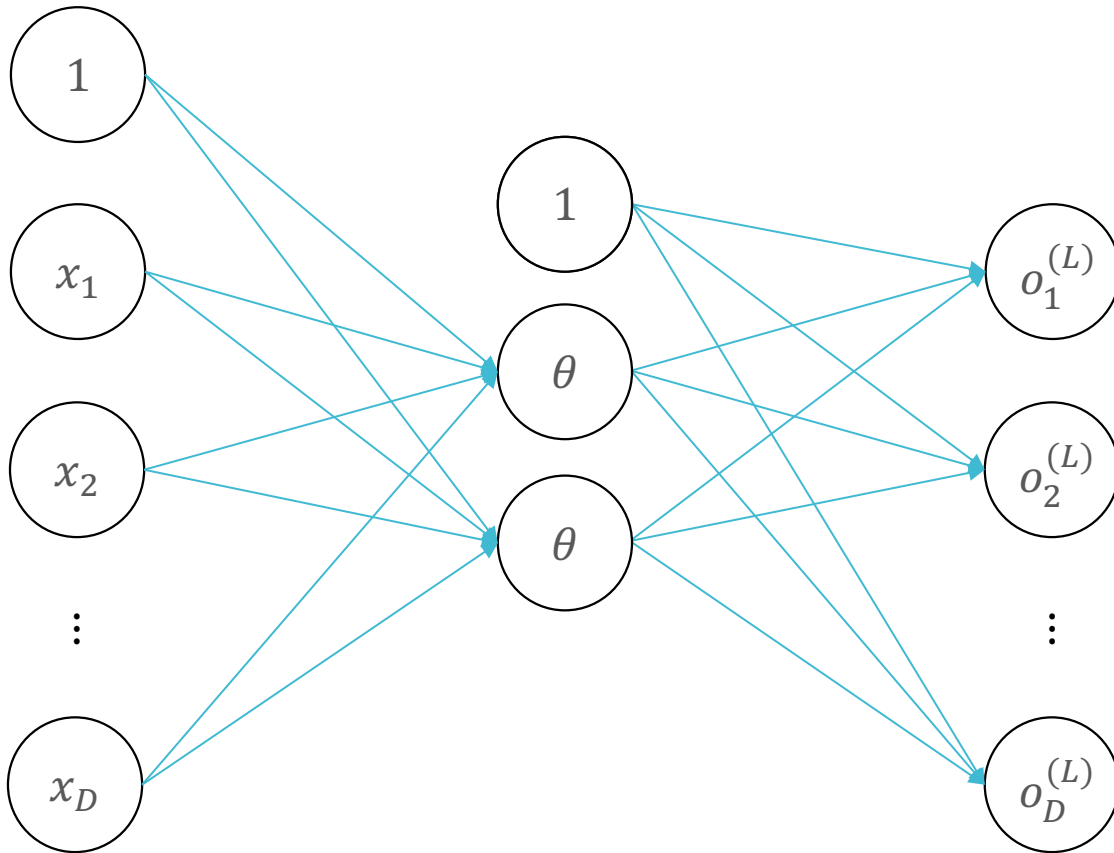
- Principal components are constrained to be orthogonal (unit) vectors

# Autoencoders



Insight: neural networks implicitly learn low-dimensional representations of inputs in hidden layers

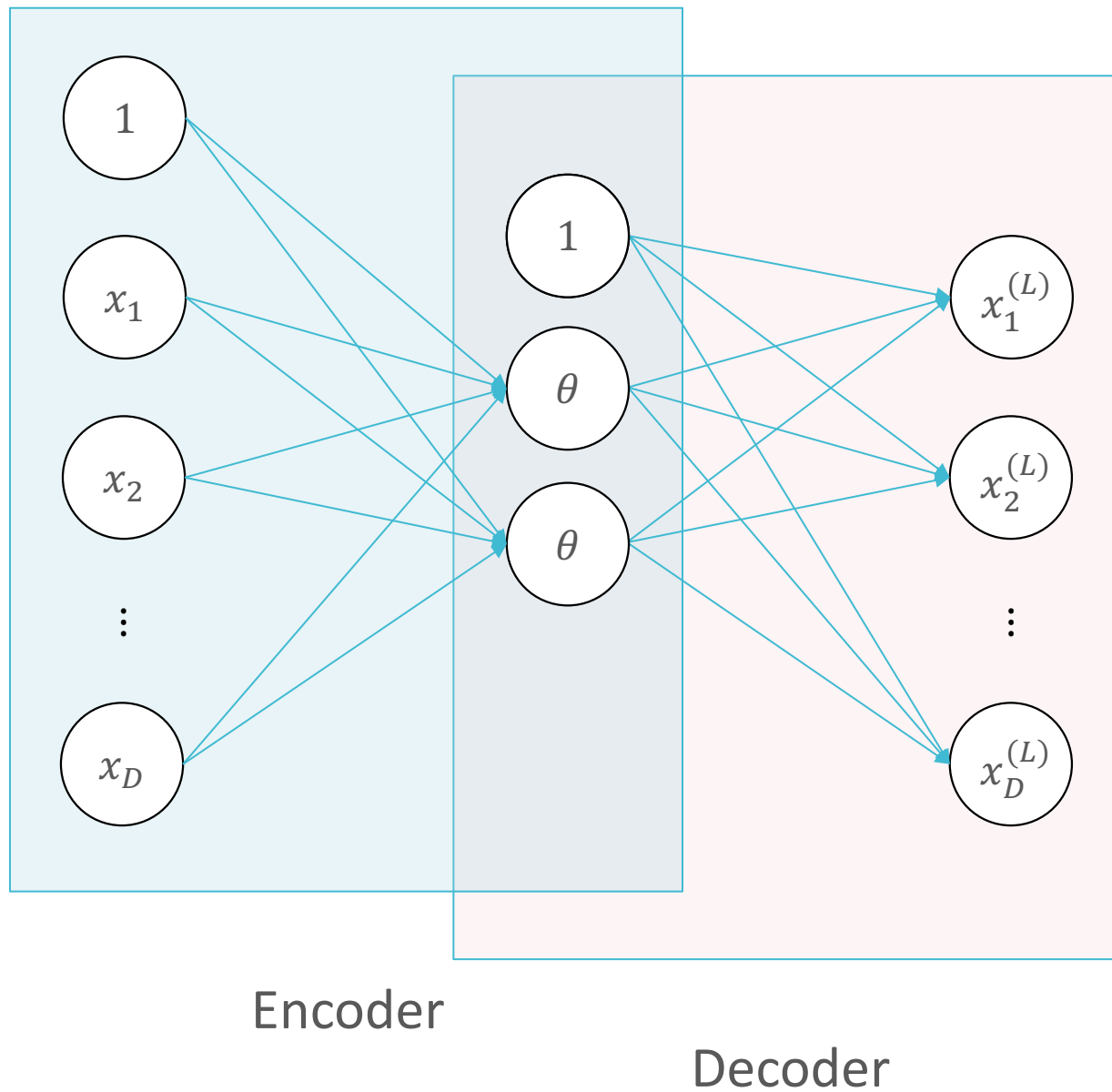
# Autoencoders



- Learn the weights by minimizing the reconstruction loss:

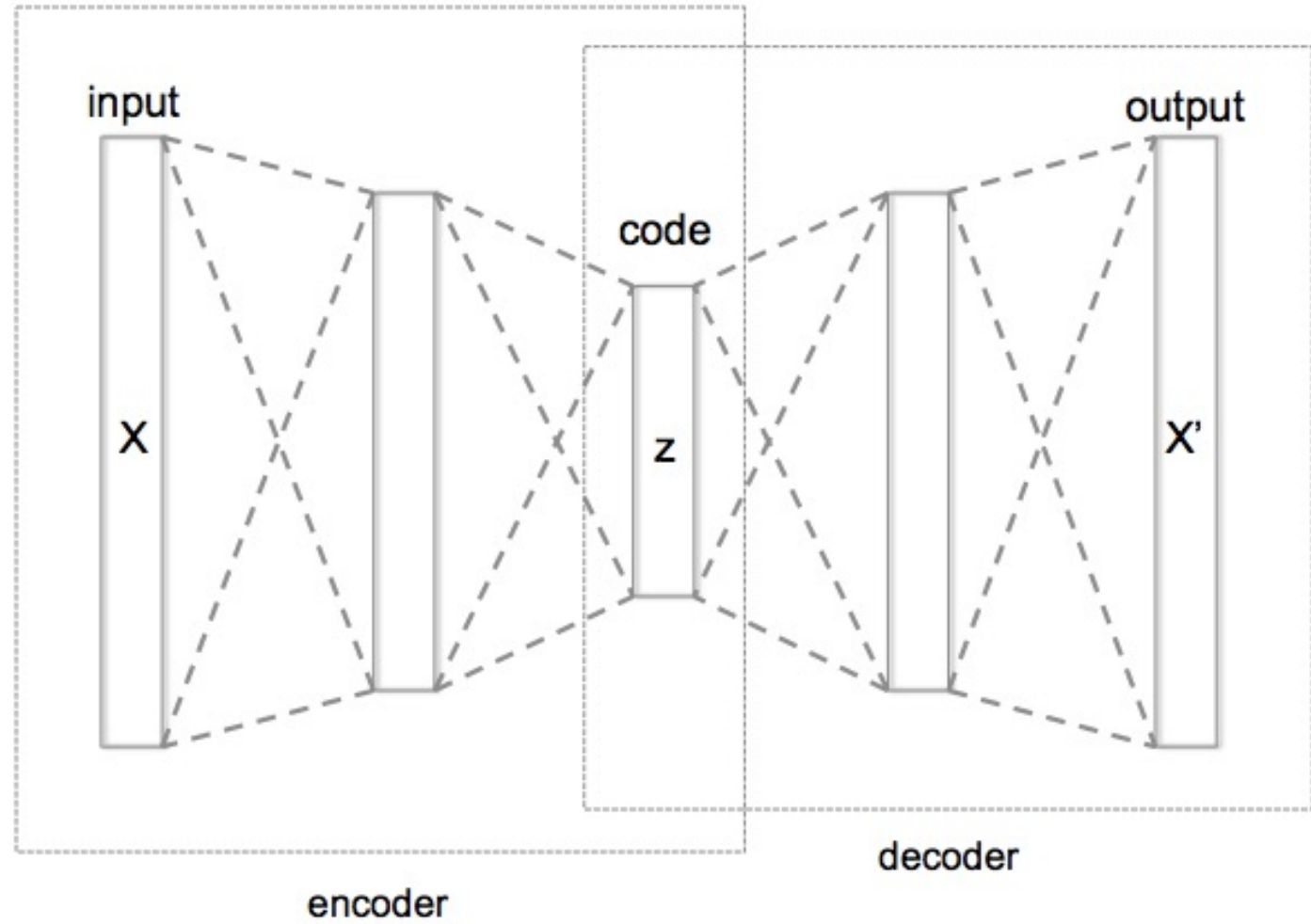
$$e(\mathbf{x}) = \|\mathbf{x} - \mathbf{o}^{(L)}\|_2^2$$

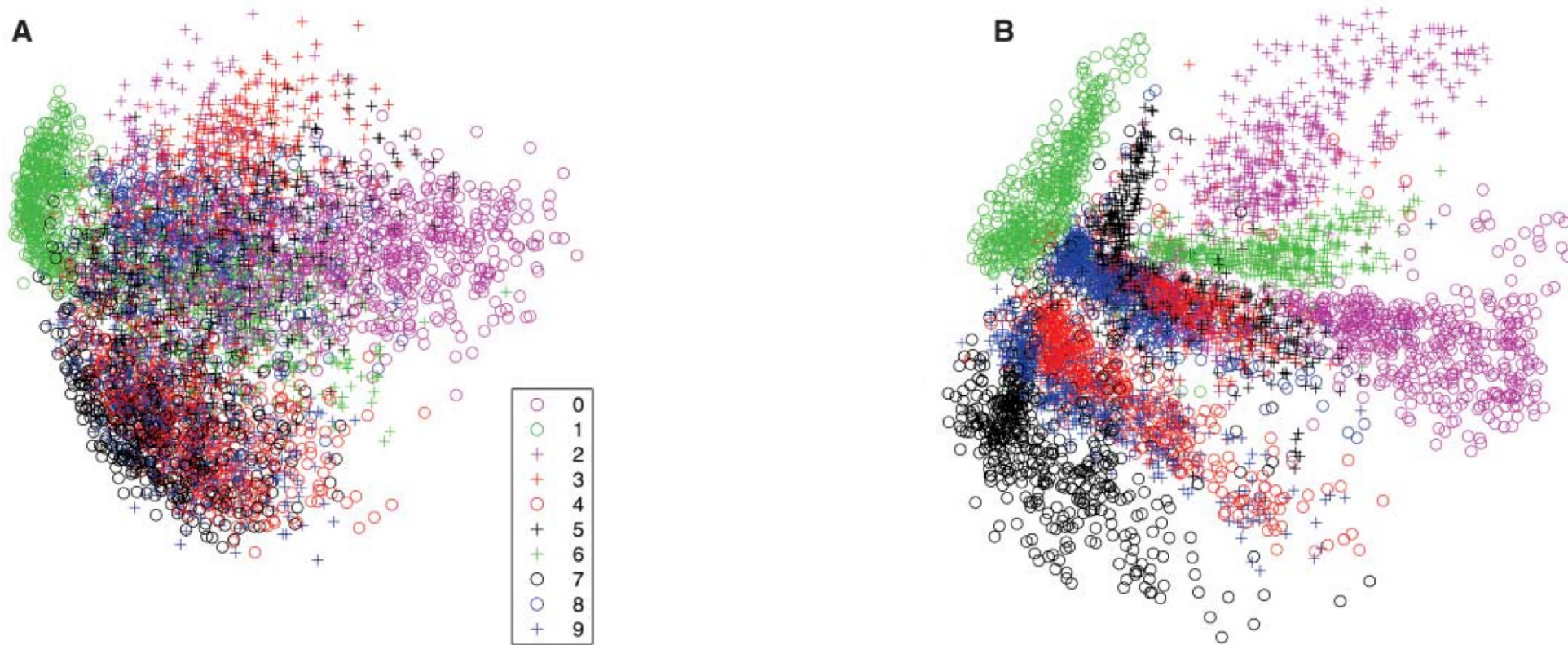
# Autoencoders





# Deep Autoencoders





# PCA (A) vs. Autoencoders (B) (Hinton and Salakhutdinov, 2014)

# Key Takeaways

- PCA finds an orthonormal basis where the first principal component maximizes the variance  $\Leftrightarrow$  minimizes the reconstruction error
  - PCs are given by the eigenvectors of the covariance matrix  $X^T X$  with the corresponding eigenvalues being a measure of the variance captured by that PC
  - Eigenvectors & eigenvalues can be computed using SVD
- ~~ICA finds statistically independent, not orthogonal components~~
- Autoencoders use neural networks to automatically learn a latent representation that minimizes the reconstruction error

# Learning Paradigms

- Supervised learning -  $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$ 
  - Regression -  $y^{(i)} \in \mathbb{R}$
  - Classification -  $y^{(i)} \in \{1, \dots, C\}$
- Unsupervised learning -  $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ 
  - Clustering
  - Dimensionality reduction
- Reinforcement learning -  $\mathcal{D} = \{(\mathbf{s}^{(n)}, \mathbf{a}^{(n)}, r^{(n)})\}_{n=1}^N$
- Active learning
- Semi-supervised learning
- Online learning

Source: <https://techobserver.net/2019/06/argo-ai-self-driving-car-research-center/>

Source: <https://www.wired.com/2012/02/high-speed-trading/>

# Reinforcement Learning: Examples



Source: <https://www.cnet.com/news/boston-dynamics-robot-dog-spot-finally-goes-on-sale-for-74500/>

Source: <https://twitter.com/alphagomovie>

# Reinforcement Learning: Problem Formulation

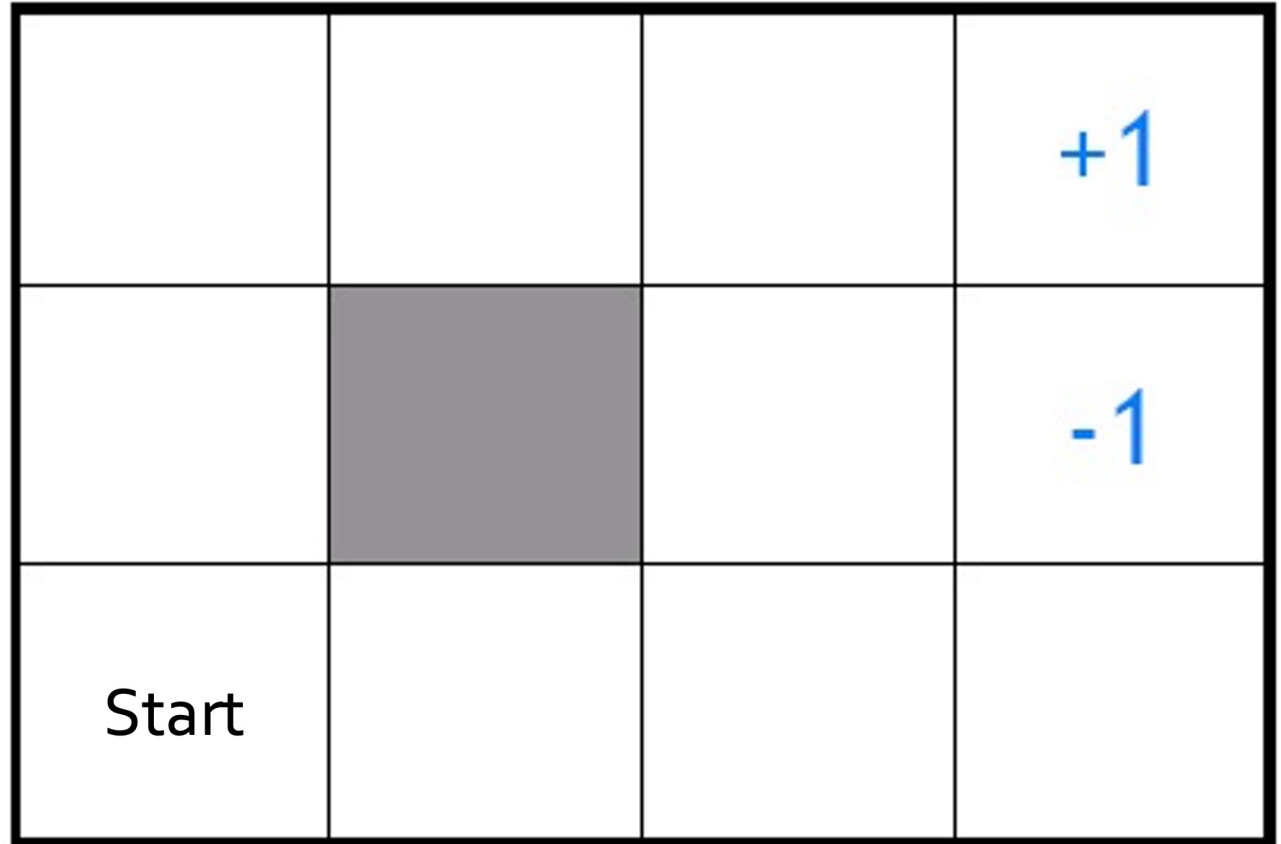
- State space,  $\mathcal{S}$
- Action space,  $\mathcal{A}$
- Reward function
  - Stochastic,  $p(r | s, a)$
  - Deterministic,  $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- Transition function
  - Stochastic,  $p(s' | s, a)$
  - Deterministic,  $\delta: \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$

# Reinforcement Learning: Problem Formulation

- Policy,  $\pi : \mathcal{S} \rightarrow \mathcal{A}$ 
  - Specifies an action to take in *every* state
- Value function,  $V^\pi : \mathcal{S} \rightarrow \mathbb{R}$ 
  - Measures the expected total payoff of starting in some state  $s$  and executing policy  $\pi$ , i.e., in every state, taking the action that  $\pi$  returns

# Toy Example

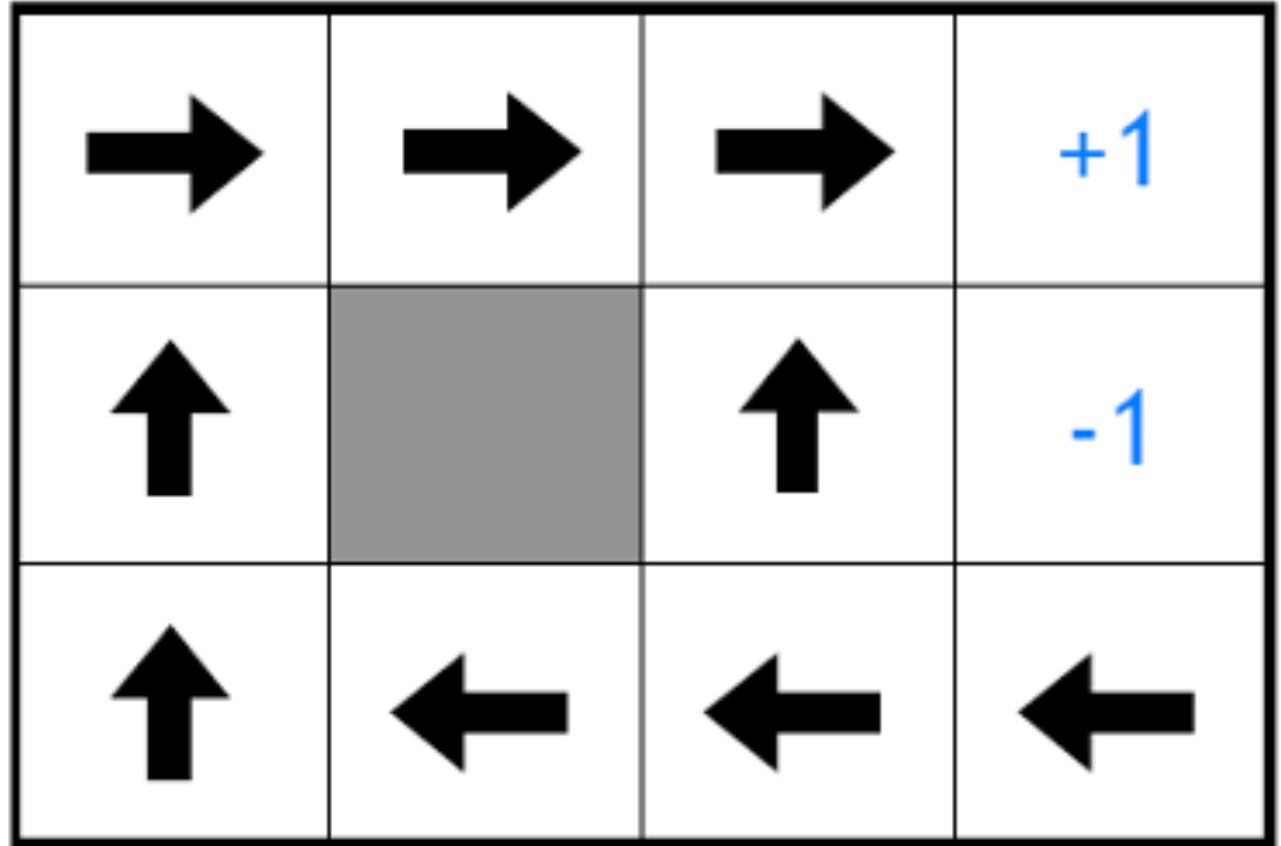
- $\mathcal{S}$  = all empty squares in the grid
- $\mathcal{A}$  = {up, down, left, right}
- Deterministic transitions
- Rewards of +1 and -1 for entering the labelled squares
- Terminate after receiving either reward





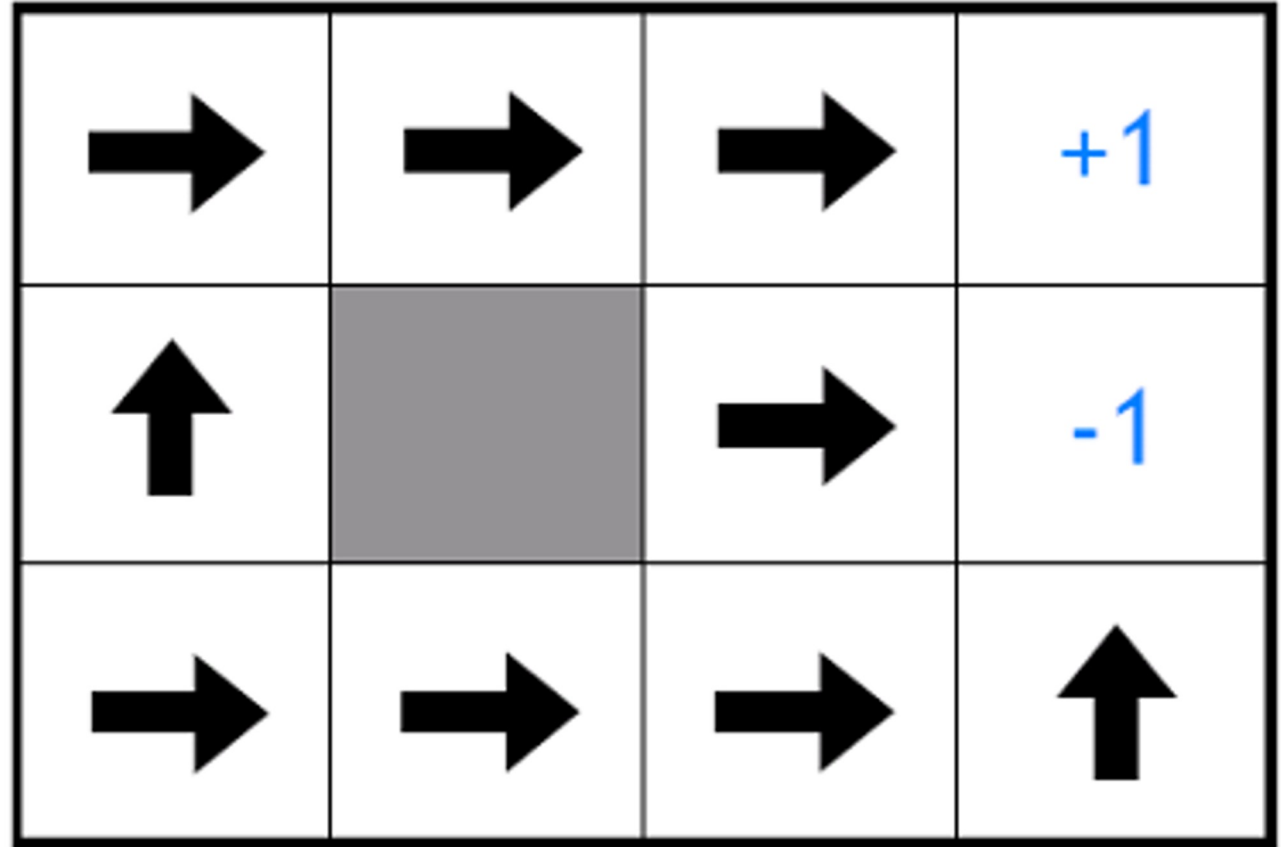
# Toy Example

Is this policy optimal?



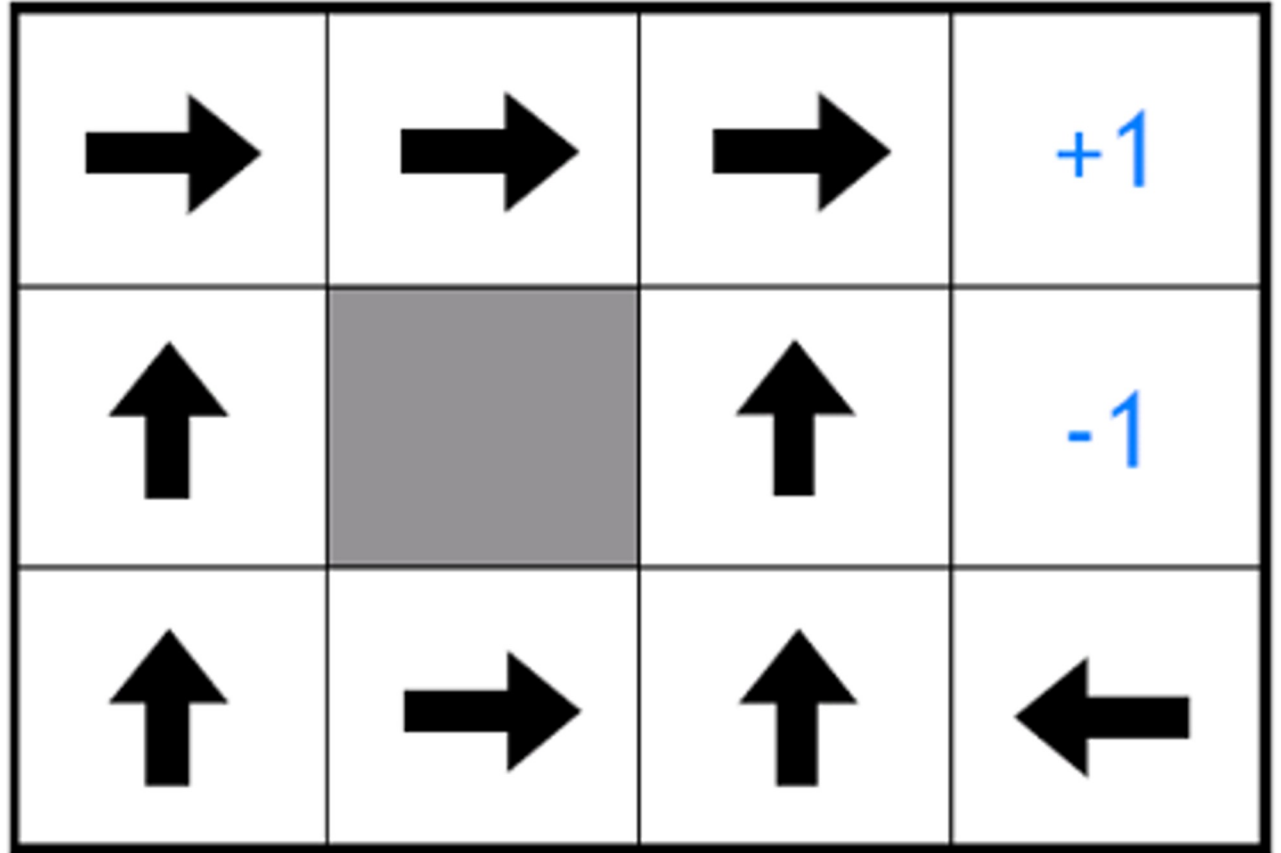
## Toy Example

Optimal policy given a  
reward of -2 per step



## Toy Example

Optimal policy given a reward of  $-0.1$  per step



# Markov Decision Process (MDP)

- Assume the following model for our data:
  1. Start in some initial state  $s_0$
  2. For time step  $t$ :
    1. Agent observes state  $s_t$
    2. Agent takes action  $a_t = \pi(s_t)$
    - 3. Agent receives reward  $r_t \sim p(r | s_t, a_t)$
    - 4. Agent transitions to state  $s_{t+1} \sim p(s' | s_t, a_t)$
  3. Total reward is  $\sum_{t=0}^{\infty} \gamma^t r_t$  discount factor  $0 \leq \gamma \leq 1$
- MDPs make the *Markov assumption*: the reward and next state only depend on the current state and action.

# Reinforcement Learning: 3 Key Challenges

1. The algorithm has to gather its own training data
2. The outcome of taking some action is often stochastic or unknown until after the fact
3. Decisions can have a delayed effect on future outcomes (exploration-exploitation tradeoff)

# MDP Example: Multi-armed bandit

- Single state:  $|\mathcal{S}| = 1$
- Three actions:  $\mathcal{A} = \{1, 2, 3\}$
- Deterministic transitions
- Rewards are stochastic

# MDP Example: Multi-armed bandit

Bandit 1	Bandit 2	Bandit 3
1	2	1
1	0	0
1	0	3
1	0	2
0	0	4
1	2	2
0	0	1
1	2	4
1	0	0
1	2	3
1	0	3
0	0	1

# Reinforcement Learning: Objective Function

- Find a policy  $\pi^* = \operatorname{argmax}_{\pi} V^{\pi}(s) \forall s \in \mathcal{S}$

•  $V^{\pi}(s) = \mathbb{E}[\text{discounted total reward of starting in state } s \text{ and executing policy } \pi \text{ forever}]$   
*(total reward of starting s ~ following  $\pi$  forever)*

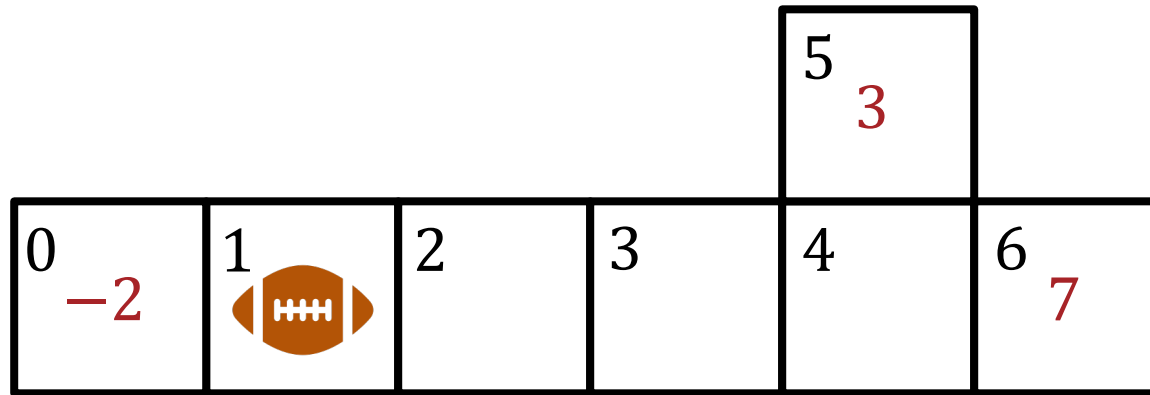
$$= \mathbb{E}_{P(s'|s,a)} [R(s_0=s, \pi(s_0)) + \gamma R(s_1, \pi(s_1)) + \gamma^2 R(s_2, \pi(s_2)) + \dots]$$

$$= \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, \pi(s_t)) \right]$$

$$= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}[R(s_t, \pi(s_t))] \leftarrow$$



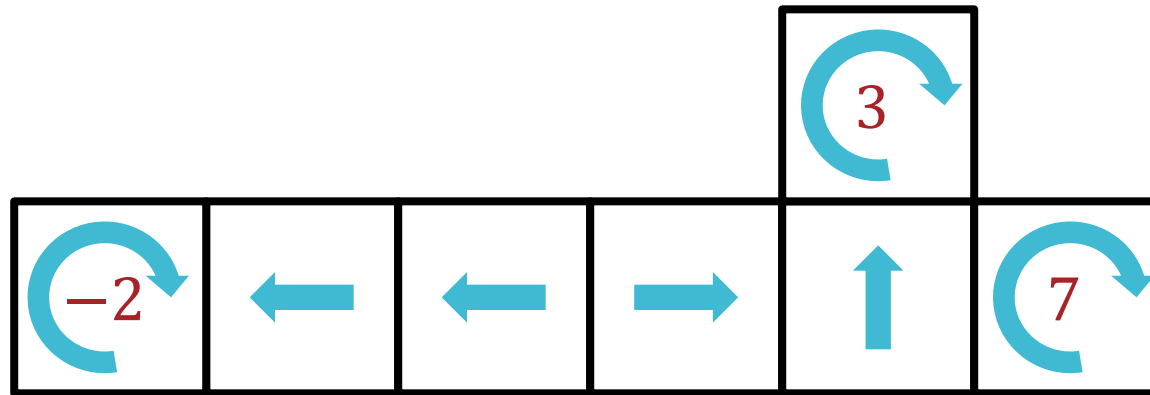
# Value Function: Example



$$R(s, a) = \begin{cases} -2 & \text{if entering state 0 (safety)} \\ 3 & \text{if entering state 5 (field goal)} \\ 7 & \text{if entering state 6 (touch down)} \\ 0 & \text{otherwise} \end{cases}$$

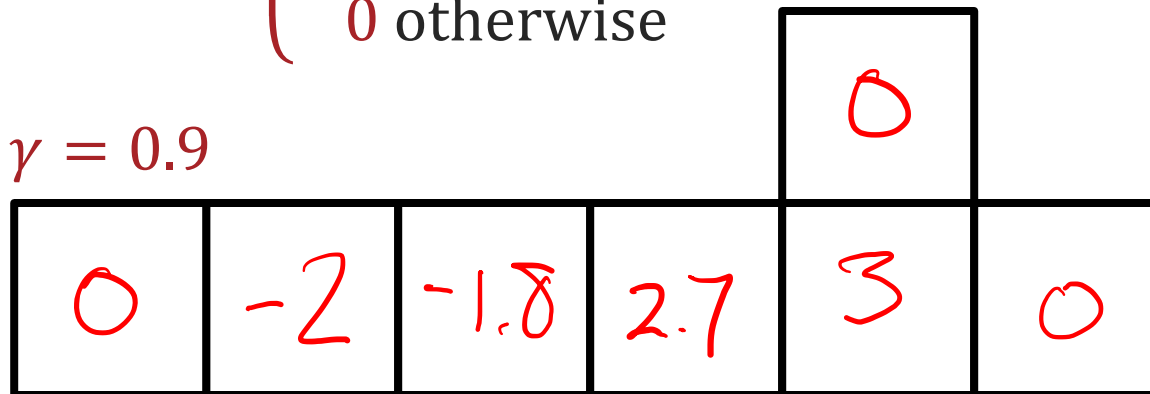
$$\gamma = 0.9$$

# Value Function: Example

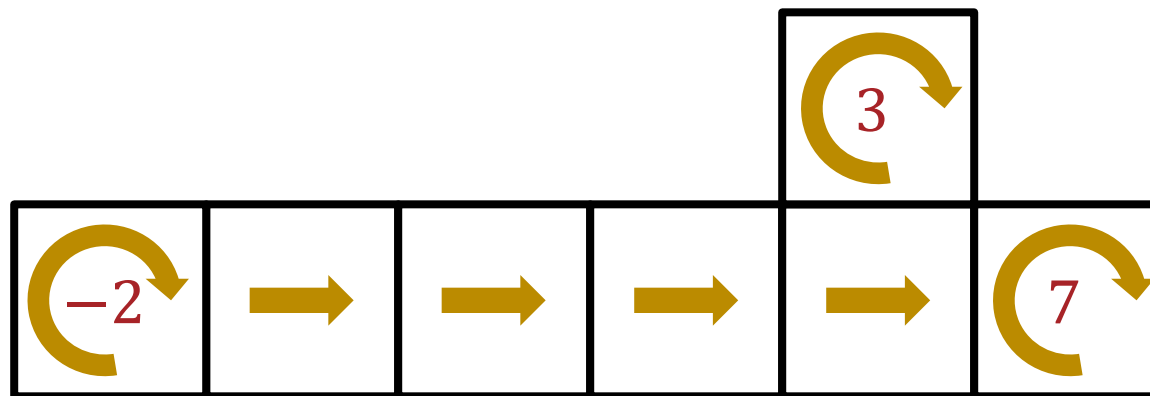


$$R(s, a) = \begin{cases} -2 & \text{if entering state 0 (safety)} \\ 3 & \text{if entering state 5 (field goal)} \\ 7 & \text{if entering state 6 (touch down)} \\ 0 & \text{otherwise} \end{cases}$$

$$\gamma = 0.9$$

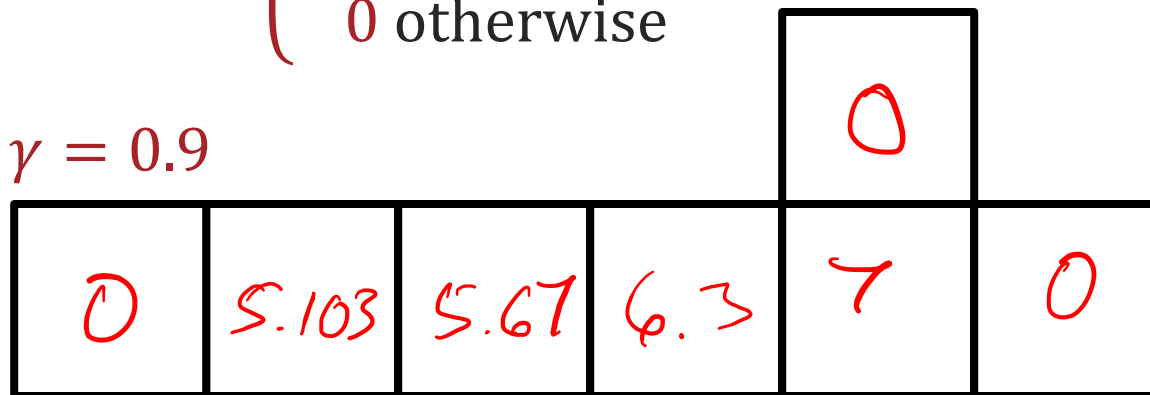


How can we  
Value Function:  
learn this  
Example  
optimal policy?



$$R(s, a) = \begin{cases} -2 & \text{if entering state 0 (safety)} \\ 3 & \text{if entering state 5 (field goal)} \\ 7 & \text{if entering state 6 (touch down)} \\ 0 & \text{otherwise} \end{cases}$$

$\gamma = 0.9$



# Value Function

- $V^\pi(s) = \mathbb{E}[\text{discounted total reward of starting in state } s \text{ and executing policy } \pi \text{ forever}]$   
 $= \mathbb{E}[R(s_0, \pi(s_0)) + \gamma R(s_1, \pi(s_1)) + \gamma^2 R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \mathbb{E}[R(s_1, \pi(s_1)) + \gamma R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \sum_{s_1 \in \mathcal{S}} p(s_1 \mid s, \pi(s)) (R(s_1, \pi(s_1)) + \gamma \mathbb{E}[R(s_2, \pi(s_2)) + \dots \mid s_1])$

# Value Function

- $V^\pi(s) = \mathbb{E}[\text{discounted total reward of starting in state } s \text{ and executing policy } \pi \text{ forever}]$   
 $= \mathbb{E}[R(s_0, \pi(s_0)) + \gamma R(s_1, \pi(s_1)) + \gamma^2 R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \mathbb{E}[R(s_1, \pi(s_1)) + \gamma R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \sum_{s_1 \in \mathcal{S}} p(s_1 \mid s, \pi(s)) (R(s_1, \pi(s_1)) + \gamma \mathbb{E}[R(s_2, \pi(s_2)) + \dots \mid s_1])$

# Value Function

- $V^\pi(s) = \mathbb{E}[\text{discounted total reward of starting in state } s \text{ and executing policy } \pi \text{ forever}]$   
 $= \mathbb{E}[R(s_0, \pi(s_0)) + \gamma R(s_1, \pi(s_1)) + \gamma^2 R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \mathbb{E}[R(s_1, \pi(s_1)) + \gamma R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \sum_{s_1 \in \mathcal{S}} p(s_1 \mid s, \pi(s)) (R(s_1, \pi(s_1)) + \gamma \mathbb{E}[R(s_2, \pi(s_2)) + \dots \mid s_1])$

# Value Function

- $V^\pi(s) = \mathbb{E}[\text{discounted total reward of starting in state } s \text{ and executing policy } \pi \text{ forever}]$   
 $= \mathbb{E}[R(s_0, \pi(s_0)) + \gamma R(s_1, \pi(s_1)) + \gamma^2 R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \mathbb{E}[R(s_1, \pi(s_1)) + \gamma R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \sum_{s_1 \in \mathcal{S}} p(s_1 \mid s, \pi(s)) (R(s_1, \pi(s_1)) + \gamma \mathbb{E}[R(s_2, \pi(s_2)) + \dots \mid s_1])$

# Value Function

- $V^\pi(s) = \mathbb{E}[\text{discounted total reward of starting in state } s \text{ and executing policy } \pi \text{ forever}]$   
 $= \mathbb{E}[R(s_0, \pi(s_0)) + \gamma R(s_1, \pi(s_1)) + \gamma^2 R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \mathbb{E}[R(s_1, \pi(s_1)) + \gamma R(s_2, \pi(s_2)) + \dots \mid s_0 = s]$   
 $= R(s, \pi(s)) + \gamma \sum_{s_1 \in \mathcal{S}} p(s_1 \mid s, \pi(s)) (R(s_1, \pi(s_1)) + \gamma \mathbb{E}[R(s_2, \pi(s_2)) + \dots \mid s_1])$

$$V^\pi(s) = R(s, \pi(s)) + \gamma \sum_{s_1 \in \mathcal{S}} p(s_1 \mid s, \pi(s)) V^\pi(s_1)$$

Bellman equations 



# Optimality

- Optimal value function:

$$V^*(s) = \max_{a \in \mathcal{A}} R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^*(s')$$

- System of  $|\mathcal{S}|$  equations and  $|\mathcal{S}|$  variables

- Optimal policy:

$$\pi^*(s) = \operatorname{argmax}_{a \in \mathcal{A}} R(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s' | s, a) V^*(s')$$

Immediate  
reward

(Discounted)  
Future reward