

## UNIT-III

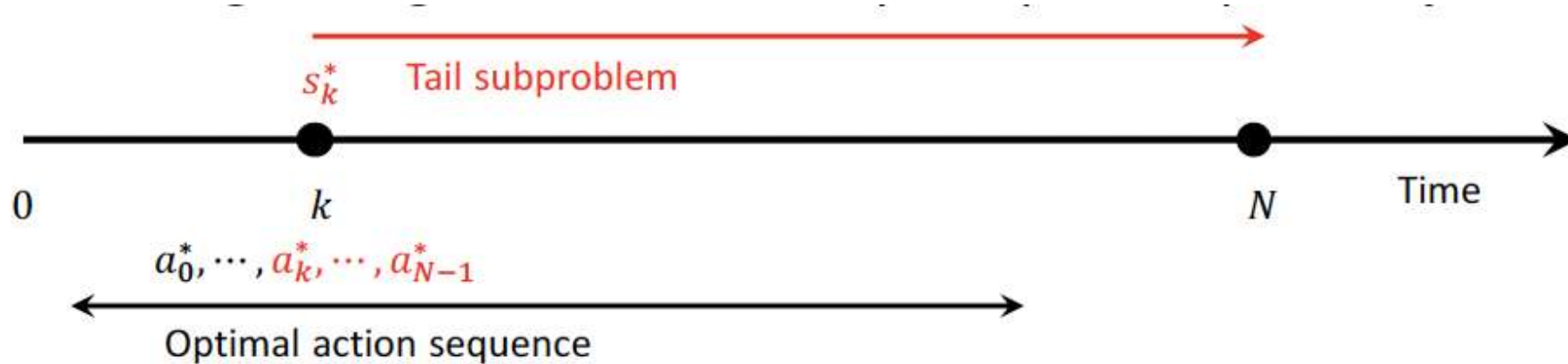
*Overview of dynamic programming for MDP, Definition and formulation of planning in MDPs, principle of optimality, Policy Evaluation, Policy Improvement, Policy Iteration, Value Iteration, Generalized Policy Iteration, Efficiency of Dynamic Programming, Banach fixed point theorem, proof of convergence of policy evaluation and value iteration algorithms*

# Dynamic Programing

- Dynamic Programing [DP] in this course, refer to a collection of algorithms that can be used to compute optimal policies given a **perfect model** of the environment in a MDP.
- Limited utility due to the 'perfect model' assumption and due to computational expense.
- But still are important as they provide essential foundation for many of the subsequent methods.
- Many of the methods can be viewed as attempts to achieve much the same effect as DP with less computation and without perfect model assumption of the environment.
- The key idea in DP is to use the value functions and Bellman equations to organize and structure the search for good policies.
- Dynamic Programing addresses a bigger problem by breaking it down as subproblems and then
  - Solving the subproblems
  - Combining solutions to subproblems

# Dynamic Programming

- *Dynamic Programming is based on the principle of optimality.*



# Requirements for Dynamic Programming

- x Optimal substructure i.e., principle of optimality applies.
- x Overlapping subproblems, i.e., subproblems recur many times and solutions to these subproblems can be cached and reused.
- x MDPs satisfy both through Bellman equations and value functions.
- x Dynamic programming is used to solve many other problems, e.g., Scheduling algorithms, Graph algorithms (e.g. shortest path algorithms), Bioinformatics etc.

# Planning by Dynamic Programming

- § Planning by dynamic programming assumes full knowledge of the MDP
- § For *prediction/evaluation*
  - ▶ Input: MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$  and policy  $\pi$
  - ▶ Output: Value function  $v_\pi$
- § For *control*
  - ▶ Input: MDP  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma \rangle$
  - ▶ Output: Optimal value function  $v_*$  and optimal policy  $\pi_*$

# Iterative Policy Evaluation

- § Problem: Policy evaluation: Compute the state-value function  $v_\pi$  for an arbitrary policy  $\pi$ .
- § Solution strategy: Iterative application of Bellman expectation equation.
- § Recall the Bellman expectation equation.

$$v_\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v_\pi(s') \right\} \quad (1)$$

- § Consider a sequence of approximate value functions  $v^{(0)}, v^{(1)}, v^{(2)}, \dots$  each mapping  $\mathcal{S}^+$  to  $\mathbb{R}$ . Each successive approximation is obtained by using eqn. (1) as an update rule.

$$v^{(k+1)}(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v^{(k)}(s') \right\}$$



# Iterative Policy Evaluation

$$v^{(k+1)}(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v^{(k)}(s') \right\}$$

- § In code, this can be implemented by using two arrays - one for the old values  $v^{(k)}(s)$  and the other for the new values  $v^{(k+1)}(s)$ . Here, new values of  $v^{(k+1)}(s)$  are computed one by one from the old values  $v^{(k)}(s)$  without changing the old values.
- § Another way is to use one array and update the values 'in place', *i.e.*, each new value immediately overwriting the old one.
- § Both these converges to the true value  $v_{\pi}$  and the 'in place' algorithm usually converges faster.

# Iterative Policy Evaluation

Iterative Policy Evaluation, for estimating  $V \approx v_\pi$

**Input:**  $\pi$ , the policy to be evaluated

**Algorithm parameter:** a small threshold  $\theta > 0$  determining accuracy of estimation

Initialize  $V(s)$ , for all  $s \in \mathcal{S}^+$ , arbitrarily except that  $V(\text{terminal}) = 0$

Loop:

$\Delta \leftarrow 0$

Loop for each  $s \in \mathcal{S}$ :

$v \leftarrow V(s)$

$$V(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) V(s') \right\}$$

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until  $\Delta < \theta$



# Evaluating a Random Policy in the Small Gridworld

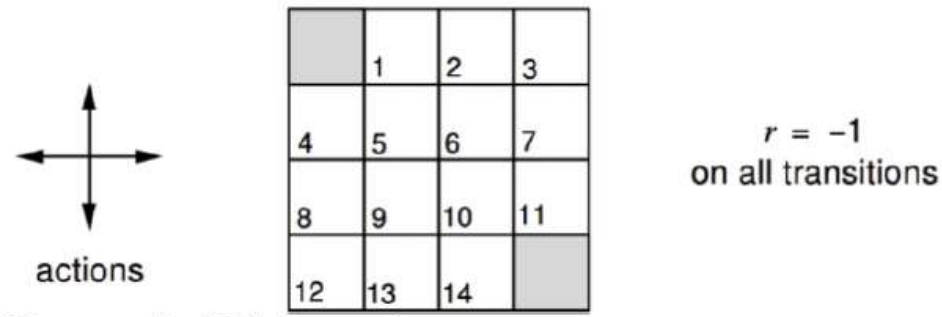


Figure credit: [SB] chapter 4

- § Undiscounted episodic MDP ( $\lambda = 1$ )
- § Non-terminal states are  $\mathcal{S} = \{1, 2, \dots, 14\}$
- § Two terminal states (shown as shaded squares)
- § 4 possible actions in each state,  $\mathcal{A} = \{up, down, right, left\}$
- § Deterministic state transitions
- § Actions leading out of the grid leave state unchanged
- § Reward is -1 until the terminal state is reached
- § Agent follows uniform random policy  
 $\pi(n|.) = \pi(s|.) = \pi(e|.) = \pi(w|.)$

We will start with initialising  $v_0$  for the random policy to all 0s.

0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0
0.0	0.0	0.0	0.0

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

$$\begin{aligned}
 v_1(6) &= \sum_{a \in \{u,d,l,r\}} \pi(a|6) \sum_{s',r} p(s',r|6,a) [r + \gamma v_0(s')] \\
 &= \sum_{a \in \{u,d,l,r\}} \underbrace{\pi(a|6)}_{= 0.25 \forall a} \sum_{s'} p(s'|6,a) \underbrace{[r + \gamma v_0(s')]}_{\substack{= -1 \\ = 0 \forall s'}} \\
 &= 0.25 * \{-p(2|6,u) - p(10|6,d) - p(5|6,l) - p(7|6,r)\} \\
 &= 0.25 * \{-1 - 1 - 1 - 1\} \\
 &= -1 \\
 &\Rightarrow v_1(6) = -1
 \end{aligned}$$

For terminal states  $p(s'/s,a) = 0$  and hence  $v_k(1) = v_k(16) = 0$  for all  $k$ . So  $v_1$  for the random policy is given by:

0.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	-1.0
-1.0	-1.0	-1.0	0.0

Now, for  $v_2(s)$  we are assuming  $\gamma$  or the discounting factor to be 1:

$$\begin{aligned}
 v_2(6) &= \sum_{a \in \{u,d,l,r\}} \underbrace{\pi(a|6)}_{= 0.25 \forall a} \sum_{s'} p(s'|6,a) \underbrace{[r + \gamma v_1(s')]}_{= -1} = \begin{cases} -1, s' \in S \\ 0, s' \in S^+ \setminus S \end{cases} \\
 &= 0.25 * \{p(2|6,u)[-1 - \gamma] + p(10|6,d)[-1 - \gamma] \\
 &\quad + p(5|6,l)[-1 - \gamma] + p(7|6,r)[-1 - \gamma]\} \\
 &\stackrel{\gamma=1}{=} 0.25 * \{-2 - 2 - 2 - 2\} \\
 &= -2
 \end{aligned}$$

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16

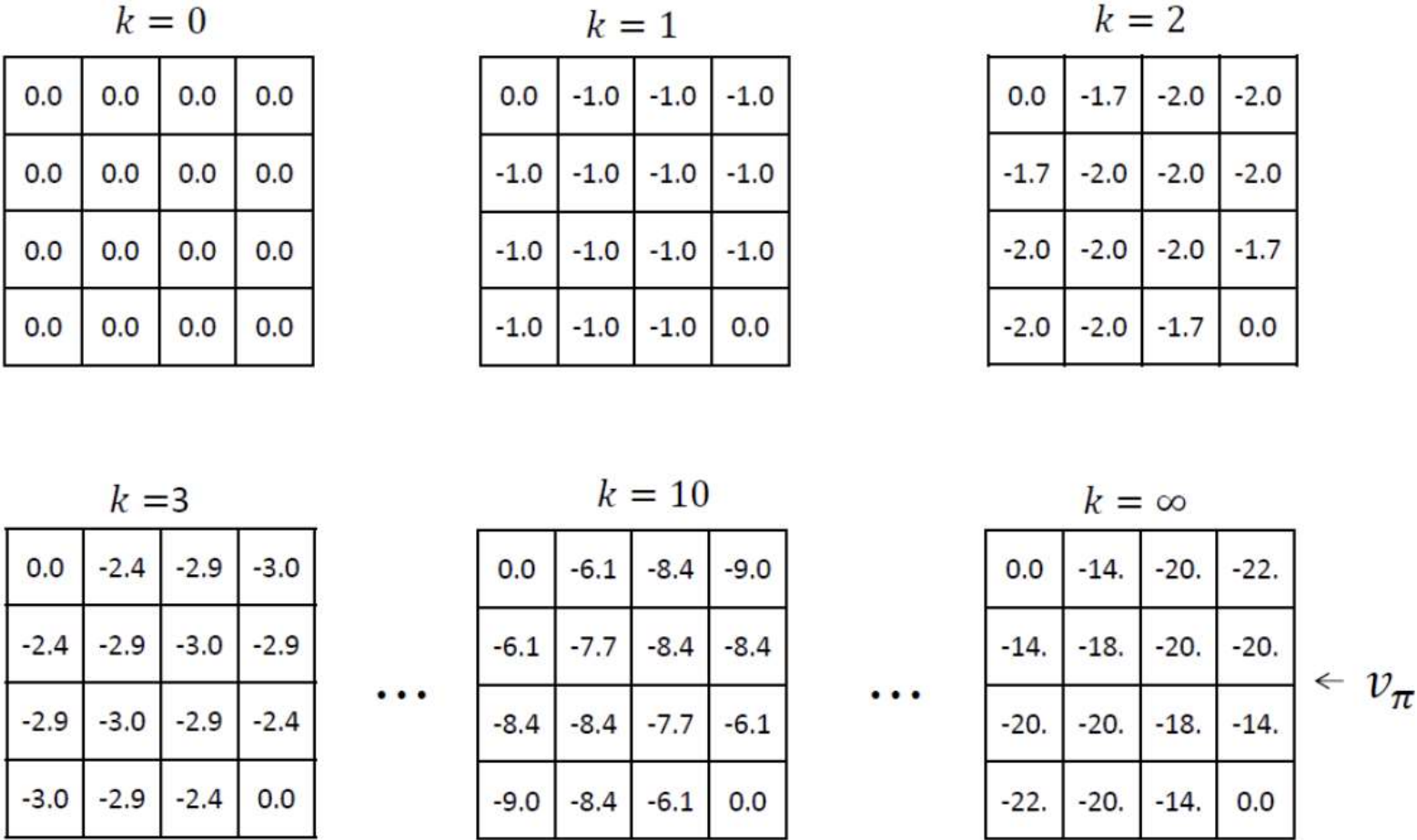
For all the remaining states, i.e., 2, 5, 12 and 15,  $v_2$  can be calculated as follows:

$$\begin{aligned}
 v_2(2) &= \sum_{a \in \{u,d,l,r\}} \underbrace{\pi(a|2)}_{= 0.25 \forall a} \sum_{s'} p(s'|2,a) [\underbrace{r}_{=-1} + \gamma \underbrace{v_1(s')}]_{=\begin{cases} -1, s' \in S \\ 0, s' \in S^+ \setminus S \end{cases}} \\
 &= 0.25 * \{p(2|2,u)[-1 - \gamma] + p(6|2,d)[-1 - \gamma] \\
 &\quad + p(1|2,l)[-1 - \gamma * 0] + p(3|2,r)[-1 - \gamma]\} \\
 &\stackrel{\gamma=1}{=} 0.25 * \{-2 - 2 - 1 - 2\} \\
 &= -1.75 \\
 &\Rightarrow v_2(2) = -1.75
 \end{aligned}$$

$\Rightarrow v_2$  for the random policy:

0.0	-1.7	-2.0	-2.0
-1.7	-2.0	-2.0	-2.0
-2.0	-2.0	-2.0	-1.7
-2.0	-2.0	-1.7	0.0

If we repeat this step several times, we get  $v_{\pi}$



# Improving a Policy: Policy Iteration

§ Given a policy  $\pi$

▶ **Evaluate** the policy

$$v_{\pi} \doteq v^{(k+1)}(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v^{(k)}(s') \right\}$$

▶ **Improve** the policy by acting greedily with respect to  $v_{\pi}$

$$\pi' = \text{greedy}(v_{\pi})$$

being greedy means choosing the action that will land the agent into best state *i.e.*,  $\pi'(s) \doteq \arg \max_{a \in \mathcal{A}} q_{\pi}(s, a) =$

$$\arg \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v_{\pi}(s') \right\}$$

§ In Small Gridworld improved policy was optimal  $\pi' = \pi_*$

§ In general, need more iterations of improvement/evaluation

§ But this process of **policy iteration** always converges to  $\pi_*$



# Policy Iteration

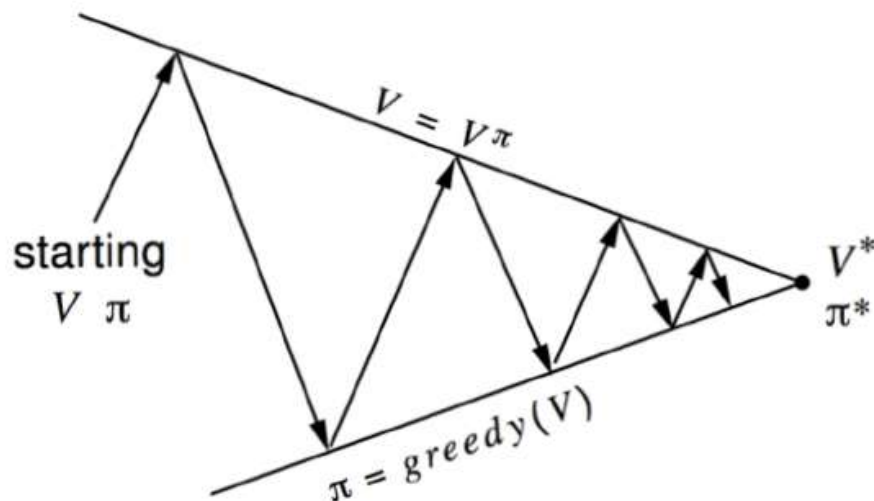
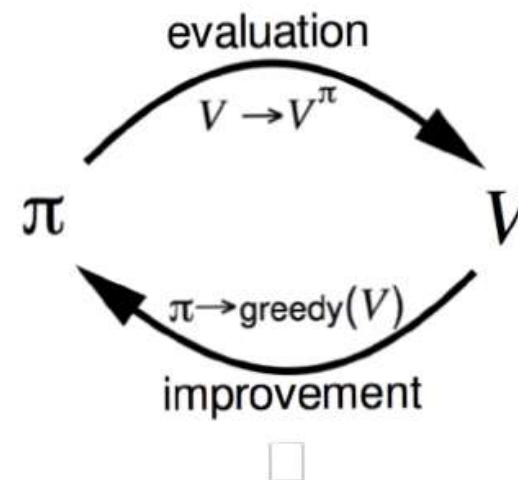


Figure credit: [David Silver: DeepMind]



§ **Policy Evaluation:** Estimate  $v_\pi$  by iterative policy evaluation.

§ **Policy Improvement:** Generate  $\pi' \geq \pi$  by greedy policy improvement.

# Policy Iteration

---

## Algorithm 1: Policy iteration

---

```
1 initialization: Select  $\pi^0, n \leftarrow 0$ ;  
2 do  
3   (Policy Evaluation)  $v_{(\pi^{n+1})} \leftarrow r_{(\pi^n)} + \gamma \mathcal{P}_{\pi^n} v_{(\pi^n)}$  ; // componentwise  
4   (Policy Improvement)  
     $\pi^{n+1}(s) \in \arg \max_{a \in \mathcal{A}} [r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v_{(\pi^{n+1})}(s')] \quad \forall s \in \mathcal{S}$ ;  
5    $n \leftarrow n + 1$ ;  
6 while  $\pi^{n+1} \neq \pi^n$ ;  
7 Declare  $\pi^* = \pi^n$ 
```

---

## Policy Iteration: Disadvantages

- § Policy iteration involves the policy evaluation step first and this itself requires a few iterations to get the exact value of  $v_\pi$  in limit.
- § The question is - must we wait for exact convergence to  $v_\pi$ ? Or can we stop short of that?
- § The small gridworld example showed that there is no change of the greedy policy after the first three iterations.
- § So the question is - is there such a number of iterations such that after that the greedy policy does not change?

# Value Iteration

- § A related question is - what about the extreme case of 1 iteration of policy evaluation and then greedy policy improvement? If we repeat this cycle, does it find the optimal policy at least in limit?
- § The good news is that - yes the guarantee is there and we will soon prove that. However, first let us modify the policy iteration algorithm to this extreme case. This is known as 'value iteration' strategy.

# Value Iteration

§ What policy iteration does: iterate over

§ And then

$$v_{\pi} \doteq v^{(k+1)}(s) \leftarrow \sum_{a \in \mathcal{A}} \pi(a|s) \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v^{(k)}(s') \right\}$$

$$\pi'(s) = \arg \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v_{\pi}(s') \right\}$$

§ What value iteration does: evaluate  $\forall a \in \mathcal{A}$

§ And then take max over it

$$r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v^{(k)}(s')$$

$$v^{(k+1)}(s) = \max_{a \in \mathcal{A}} \left\{ r(s, a) + \gamma \sum_{s' \in \mathcal{S}} p(s'|s, a) v^{(k)}(s') \right\} \text{ Where have we seen it?}$$



# Value Iteration

---

**Algorithm 2:** Value iteration

---

```
8 initialization:  $v \leftarrow v^0 \in \mathcal{V}$ , pick an  $\epsilon > 0$ ,  $n \leftarrow 0$ ;  
9 while  $\|v^{n+1} - v^n\| > \epsilon \frac{1-\gamma}{2\gamma}$  do  
10   foreach  $s \in \mathcal{S}$  do  
11      $v^{n+1}(s) \leftarrow \max_a \left\{ r(s, a) + \gamma \sum_{s'} p(s'/s, a) v^n(s') \right\}$   
12   end  
13    $n \leftarrow n + 1$ ;  
14 end  
15 foreach  $s \in \mathcal{S}$  do  
16   /* Note the use of  $\pi(s)$ . It means deterministic policy */  
    $\pi(s) \leftarrow \arg \max_a \left\{ r(s, a) + \gamma \sum_{s'} p(s'/s, a) v^n(s') \right\}$  ;    //  $n$  has  
   already been incremented by 1  
17 end
```

---



# Summary of Exact DP Algorithms for Planning

<b>Problem</b>	<b>Bellman Equation</b>	<b>Algorithm</b>
Prediction	Bellman Expectation Equation	Iterative Policy Evaluation
Control	Bellman Expectation Equation + Greedy Policy Improvement	Policy Iteration
Control	Bellman Optimality Equation	Value Iteration

## Definition

Given a vector space  $\mathcal{V} \subseteq \mathbb{R}^d$ , a function  $f : \mathcal{V} \rightarrow \mathbb{R}^+$  is a norm (denoted as  $||\cdot||$ ) if and only if

§  $||v|| \geq 0 \quad \forall v \in \mathcal{V}$

§  $||v|| = 0$  if and only if  $v = 0$

§  $||\alpha v|| = |\alpha| ||v||, \forall \alpha \in \mathbb{R} \text{ and } \forall v \in \mathcal{V}$

§ **Triangle inequality:**  $||u + v|| \leq ||u|| + ||v|| \quad u, v \in \mathcal{V}$

# Different types of Norms

§  $L_p$  norm:

$$\|v\|_p = \left( \sum_{i=1}^d |v_i|^p \right)^{\frac{1}{p}}$$

§  $L_0$  norm:

$$\|v\|_0 = \text{Number of non-zero elements in } v$$

§  $L_\infty$  norm:

$$\|v\|_\infty = \max_{1 \leq i \leq d} |v_i|$$

# Cauchy Sequence, Completeness

## Definition

A sequence of vectors  $v_1, v_2, v_3, \dots \in \mathcal{V}$  (with subscripts  $n \in \mathbb{N}$ ) is called a **Cauchy sequence** if for any positive real  $\epsilon > 0, \exists N \in \mathbb{Z}^+$  such that  $\forall m, n > N, \|v_m - v_n\| < \epsilon$ .

- § Basically, for any real positive  $\epsilon$ , an element can be found in the sequence, beyond which any two elements of the sequence will be within  $\epsilon$  of each other.
- § In other words, the elements of the sequence comes closer and closer to each other - *i.e.*, the sequence converges.

## Definition

A vector space  $\mathcal{V}$  equipped with a norm  $\|\cdot\|$  is **complete** if **every** Cauchy sequence converges in that norm to a point **in the space**. To pay tribute to Stefan Banach, the great Polish mathematician, such a space is also called the **Banach space**.

# Contraction Mapping, Fixed Point

## Definition

An operator  $\mathcal{T} : \mathcal{V} \rightarrow \mathcal{V}$  is **L-Lipschitz** if for any  $u, v \in \mathcal{V}$

$$\|\mathcal{T}u - \mathcal{T}v\| \leq L\|u - v\|$$

§ If  $L \leq 1$ , then  $\mathcal{T}$  is called a **non-expansion**, while if  $0 \leq L < 1$ , then  $\mathcal{T}$  is called a **contraction**.

## Definition

Let  $v$  is a vector in the vector space  $\mathcal{V}$  and  $\mathcal{T}$  is an operator  $\mathcal{T} : \mathcal{V} \rightarrow \mathcal{V}$ . Then  $v$  is called a **fixed point** of the operator  $\mathcal{T}$ , if  $\mathcal{T}v = v$ .



# Banach Fixed Point Theorem

## Theorem

Suppose  $\mathcal{V}$  is a Banach space and  $T : \mathcal{V} \rightarrow \mathcal{V}$  is a contraction mapping, then,

- $\exists$  an unique  $v^*$  in  $\mathcal{V}$  s.t.  $Tv^* = v^*$  and
- for arbitrary  $v^0$  in  $\mathcal{V}$ , the sequence  $\{v^n\}$  defined by  $v^{n+1} = Tv^n = T^{n+1}v^0$ , converges to  $v^*$ .

The above theorem tells that

- §  $T$  has fixed point, **an unique fixed point**.
- § For **arbitrary** starting point if we keep repeatedly applying  $T$  on that starting point, then we will converge to  $v^*$ .



# Banach Fixed Point Theorem - Proof (1)

§ Let  $v^n$  and  $v^{m+n}$  be two values of  $v$  obtained after the  $n^{th}$  and the  $(n + m)^{th}$  iteration.

$$\begin{aligned} \|v^{m+n} - v^n\| &\leq \sum_{k=0}^{m-1} \|v^{n+k+1} - v^{n+k}\| \quad [\text{Triangle inequality}] \\ &= \sum_{k=0}^{m-1} \|T^{n+k}v^1 - T^{n+k}v^0\| \leq \sum_{k=0}^{m-1} \lambda \|T^{n+k-1}v^1 - T^{n+k-1}v^0\| \\ &\leq \sum_{k=0}^{m-1} \lambda^{n+k} \|v^1 - v^0\| \quad [\text{Repeated use of contraction}] \\ &= \|v^1 - v^0\| \sum_{k=0}^{m-1} \lambda^{n+k} \\ &= \frac{\lambda^n(1 - \lambda^m)}{1 - \lambda} \|v^1 - v^0\| \end{aligned} \tag{3}$$

## Banach Fixed Point Theorem - Proof (2)

- § As  $m$  and  $n \rightarrow \infty$  and as  $\lambda < 1$ , the norm of difference between  $v^{m+n}$  and  $v^n$  becomes less and less.
- § That means the sequence  $\{v^n\}$  is Cauchy.
- § And since  $\mathcal{V}$  is a Banach space and since every Cauchy sequence converges to a point in that Banach space, therefore the Cauchy sequence  $\{v^n\}$  also converges to a point in  $\mathcal{V}$ .
- § What we have proved till now is that the sequence  $\{v^n\}$  will reach a converging point in the same space.
- § Lets say that the converging point is  $v^*$ .
- § What we will try to prove next is that  $v^*$  is a fixed point and then we will try to prove that  $v^*$  is an **unique** fixed point.

# Banach Fixed Point Theorem - Proof (3)

- § Let us try to see what we get as the norm of the difference between  $v^*$  and  $Tv^*$ .
- § In the first line below we apply triangle inequality where  $v^n$  is the value of  $v$  at the  $n^{th}$  iteration.

$$\begin{aligned} \|Tv^* - v^*\| &\leq \|Tv^* - v^n\| + \|v^n - v^*\| \\ &= \|Tv^* - Tv^{n-1}\| + \|v^n - v^*\| \\ &\leq \lambda \|v^* - v^{n-1}\| + \|v^n - v^*\| \quad [\text{Contraction property}] \end{aligned} \tag{4}$$

- § Since  $\{v^n\}$  is Cauchy and  $v^*$  is the convergence point, both the terms in the above equation will tend to 0 as  $n \rightarrow \infty$ .
- § So, as  $n \rightarrow \infty$ ,  $\|Tv^* - v^*\| \rightarrow 0$ . That means in limit  $Tv^* = v^*$ . So, it is proved that  $v^*$  is a fixed point.

## Banach Fixed Point Theorem - Proof (4)

- § Now we will show the uniqueness, *i.e.*,  $v^*$  is unique.
- § Let  $u^*$  and  $v^*$  be two fixed points of the space. From the contraction property, we can write  $\|Tu^* - Tv^*\| \leq \lambda \|u^* - v^*\|$ .
- § But, since  $u^*$  and  $v^*$  are fixed points,  $Tu^* = u^*$  and  $Tv^* = v^*$ .
- § That means  $\|u^* - v^*\| \leq \lambda \|u^* - v^*\|$  which can not be true for  $\lambda < 1$  unless  $v^* = u^*$ .
- § So, it is proved that  $v^*$  is an unique fixed point.