The idea of Policy Iteration algorithm can be applied using both Dynamic Programming and Monte Carlo Simulation, although both can be proved to have policy improvement at each iteration, the algorithm that has such guarantee is not exactly the same, so are the proofs.

The document is a comparison of the two algorithms used and their corresponding proof, in the end we show why the PI algorithm used in the dynamic programming case would fail in the Monte Carlo case.

Note that Alg 2 has multiple variants, including non-stationary $\epsilon$(an example is **GLIE**:Greedy in the Limit with Infinite Exploration), incremental implementation, fix step size $\alpha$, first-visit/every-visit MC etc. Here we only show one of them.

1. The basic policy iteration algorithm when given the MDP is as Alg 1

---

**Algorithm 1:** Policy Iteration

**Input:** Markov Decision Process $\langle S, A, P, \mathcal{R}, \gamma \rangle$
**Output:** $V^*, \pi^*$
    *Initialisation* : $V(s) \in \mathbb{R}, \pi(s) \in A(s)$. arbitrarily, $\forall s \in S$
1: *policy-stable* $\leftarrow$ `False`
2: **while** *policy-stable* $=$ `False` **do**

3:    Iterative **Policy Evaluation**: $V(s) \leftarrow V^\pi(s)$

4:    $\Delta \leftarrow \theta$
5:    **while** $\Delta \geq \theta$ **do**
6:      **for** $s \in S$ **do**
7:        $v \leftarrow V(s)$
8:        $V(s) \leftarrow \sum_{a \in A} \pi(a \mid s) \left( R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V(s') \right)$
9:        $\Delta \leftarrow \max\left( \Delta, \mid v - V(s) \mid \right)$
10:    **end for**
11:    **end while**

12:    **Policy Improvement**: $\pi \leftarrow$ `greedy`$(V)$

13:    *policy-stable* $\leftarrow$ `False`
14:    **for** $s \in S$ **do**
15:      *previous-action* $\leftarrow \pi(s)$
16:      $\pi(s) \leftarrow \text{argmax}_a \sum R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V(s')$
17:      **if** *previous-action* $\neq \pi(s)$ **then**
18:        *policy-stable* $\leftarrow$ `False`
19:      **end if**
20:    **end for**
21: **end while**
22: **return** $V, \pi$

---

The proof of policy improvement is as follows:

*Proof.* Denote the previous policy as $\pi$ and the updated policy as $\pi'$. Note that according to Algorithm 1, in every policy improvement step, we have $\pi'(s) \leftarrow \text{argmax}_a q_\pi(s, a)$.

$$
\begin{aligned}
v_\pi(s) &= q_\pi(s, \pi(s)) \\
&\leq q_\pi(s, \pi'(s)) \\
&= \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma v_\pi(S_{t+1}) \mid S_t = s\right] \\
&\leq \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma q_\pi(S_{t+1}, \pi'(S_{t+1})) \mid s_t = s\right] \\
&\vdots \\
&\leq \mathbb{E}_{\pi'}\left[R_{t+1} + \gamma R_{t+2} + \cdots \mid S_t = s\right] \\
&= v_{\pi'}(s)
\end{aligned}
$$

∎

2. The generalized policy iteration algorithm used for Monte-Carlo control is as Alg 1. Note that this is an **on-policy** algorithm which attempts to evaluate/improve the policy that's used to make decisions.

   It's important to note the two differences of this algorithm compared with Algorithm 1:

   - we're now using action-value function instead of state-value function. This is because we don't get access to the dynamics(transition and reward) of the MDP anymore. The Bellman Optimality Equation for $V(s)$ requires $R$ and $P$ to perform the one step "look-ahead".

   - we're now using $\epsilon-$greedy policy improvement instead of greedy. This is for the consideration of exploration. Otherwise, if at state $s$ there exists a action $a_1$ that in expectation has high reward but when we first visit the state we get 0 reward, at the same time there's another action $a_2$ which always gives us positive reward 1, the greedy algorithm will thus always select $a_2$.

   First we state the $\epsilon-$greedy policy improvement theorem:

   **Theorem 1.** *For any $\epsilon$-greedy policy $\pi$, the $\epsilon$-greedy policy $\pi'$ with respect to $q_\pi$ is an improvement, $v_{\pi'}(s) \geq v_\pi(s)$.*

   The proof is as follows:

   *Proof.* Denote the previous policy as $\pi$ and the updated policy as $\pi'$. Note that $\pi$ can be any $\epsilon-$greedy policy. And we observe the following equality:

   $$
   \sum_a \frac{\pi(a \mid s) - \frac{\epsilon}{|A|}}{1 - \epsilon} = \frac{1}{1 - \epsilon} - \frac{\epsilon}{1 - \epsilon} = 1
   $$

---

**Algorithm 2:** Generalized Policy Iteration for On-policy Every-visit Monte-Carlo Control

---

**Input:** State space $S$, action space $A$, environment interface $(s', r) \leftarrow f(s, a)$, discount factor $\gamma$, exploration probability $\epsilon$

**Output:** $\pi^*$

    *Initialisation* : $Q(s, a) \in \mathbb{R}, \forall s \in S, a \in A.$ arbitrarily, $N_{episode} \leftarrow 0$

1: **while** want new episode **do**
2:      $N_episode \leftarrow N_episode + 1$
3:      Sample an episode following $\pi$: $\{S_1, A_1, R_1, S_2, \cdots, S_{T-1}, A_{T-1}, R_{T-1}, S_T\}$

4:      Monte-Carlo **Policy Evaluation** with action-value function:

5:      **for** $t \in [1, \cdots, T)$ **do**
6:        $N_{(S_t, A_t)} \leftarrow N_{(S_t, A_t)} + 1$
7:        $G_t \leftarrow \sum_t^{T-1} R_t$
8:        $Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \frac{G_t - Q(S_t, A_t)}{N_{(S_t, A_t)}}$
9:      **end for**

10:     $\epsilon-$greedy **Policy Improvement**:

11:     **for** $s \in S$ **do**
12:       $A^* \leftarrow \operatorname{argmax}_a Q(s, a)$
13:       **for** $a \in A$ **do**
14:         **if** $a = A^*$ **then**
15:           $\pi(a|s) \leftarrow 1 - \epsilon + \frac{\epsilon}{|A|}$
16:         **else**
17:           $\pi(a|s) \leftarrow \frac{\epsilon}{|A|}$
18:         **end if**
19:       **end for**
20:     **end for**
21: **end while**
22: **return** $\pi$

---

$$q_\pi(s, \pi'(s)) = \sum_a \pi'(a \mid s) q_\pi(s, a)$$

$$= \frac{\epsilon}{\mid A \mid} \sum_a q_\pi(s, a) + (1 - \epsilon) \max_a q_\pi(s, a)$$

$$\geq \frac{\epsilon}{\mid A \mid} \sum_a q_\pi(s, a) + \sum_a \frac{\pi(a \mid s) - \frac{\epsilon}{|A|}}{1 - \epsilon} q_\pi(s, a)$$

$$= \sum_a \pi(a \mid s) q_\pi(s, a)$$

$$= v_\pi(s)$$

With the above inequality, we simply use the same trick to get $v_{\pi'}(s) \geq q_\pi(s, \pi'(s))$ and thus $v_{\pi'}(s) \geq v_\pi(s)$.

$\blacksquare$