$$V^\pi(s) \quad = \quad \sum_{t=1}^{\infty} E\left\{\gamma^{t-1} r_t \mid s_0 = s\right\} \tag{1}$$

$$= \quad \sum_a \pi(s,a) \left[ R(s,a) + \gamma \sum_{s'} P(s,s',a) V^\pi(s') \right] \tag{2}$$

$$d^\pi(s') \quad = \quad \lim_{t\to\infty} Pr\left\{s_t = s' \mid s_0, \pi\right\} \qquad \text{(does not depend on } s_0\text{)} \tag{3}$$

$$= \quad \sum_s d^\pi(s) \sum_a \pi(s,a) P(s,s',a) \tag{4}$$

$$\rho^\pi \quad = \quad \lim_{T\to\infty} \frac{1}{T} \sum_{t=1}^{T} r_t \qquad \text{(does not depend on } s_0\text{)} \tag{5}$$

$$= \quad \sum_s d^\pi(s) \sum_a \pi(s,a) R(s,a) \tag{6}$$

In trying to form an overall discounted performance measure for $\pi$, can we use $J(\pi) = \sum_s d^\pi(s) V^\pi(s)$? It turns out we then end up with no effect of the discounting:

$$J(\pi) \quad = \quad \sum_s d^\pi(s) V^\pi(s) \tag{7}$$

$$= \quad \sum_s d^\pi(s) \sum_a \pi(s,a) \left[ R(s,a) + \gamma \sum_{s'} P(s,s',a) V^\pi(s') \right] \tag{8}$$

$$= \quad \rho^\pi + \gamma \sum_s d^\pi(s) \sum_a \pi(s,a) \sum_{s'} P(s,s',a) V^\pi(s') \tag{9}$$

$$= \quad \rho^\pi + \gamma \sum_{s'} V^\pi(s') \sum_s d^\pi(s) \sum_a \pi(s,a) P(s,s',a) \tag{10}$$

$$= \quad \rho^\pi + \gamma \sum_{s'} V^\pi(s') d^\pi(s') \tag{11}$$

$$= \quad \rho^\pi + \gamma J(\pi) \tag{12}$$

$$= \quad \rho^\pi + \gamma \rho^\pi + \gamma^2 J(\pi) \tag{13}$$

$$= \quad \rho^\pi + \gamma \rho^\pi + \gamma^2 \rho^\pi + \gamma^3 \rho^\pi + \cdots \tag{14}$$

$$= \quad \frac{1}{1-\gamma} \rho^\pi \tag{15}$$

which is basically a scaled $\rho^\pi$, with no effect of discounting.