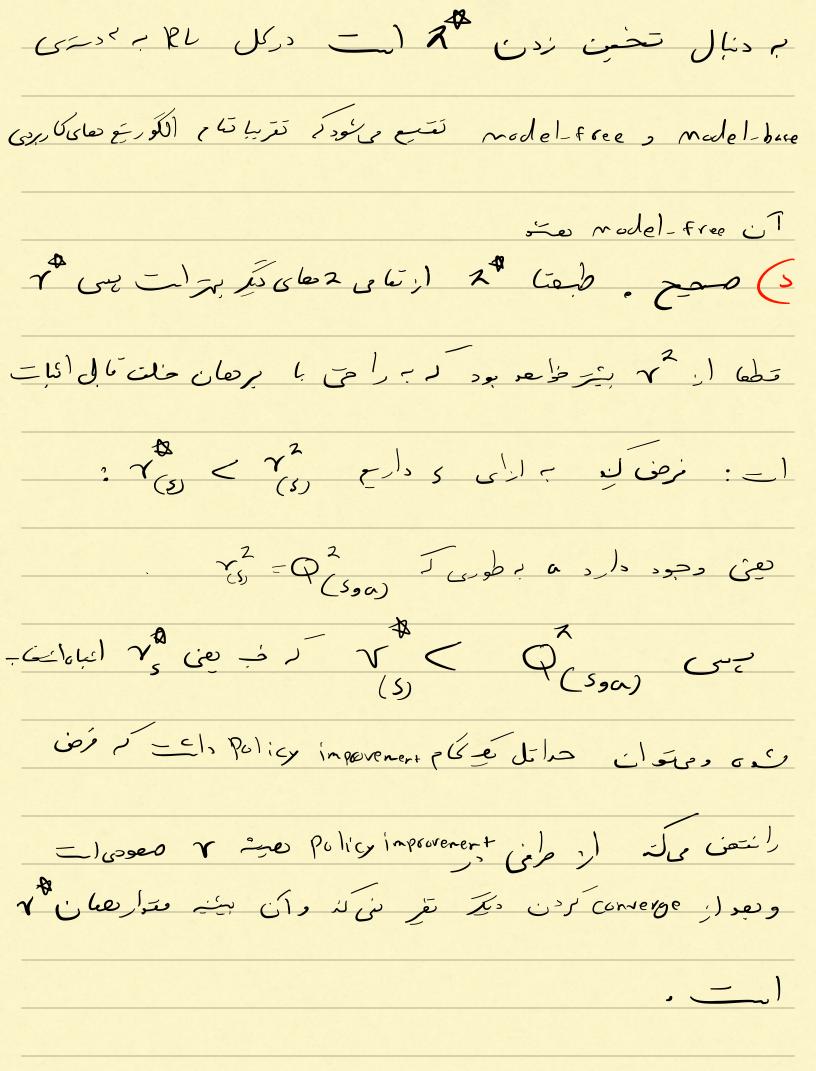
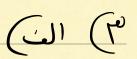
تَوِين سوم هوسي حد بام تا بنی ۱۹۸۶۸ و ۵۵۱ ا) الصحیح اے در RJ ، Tc R رانی دانے و در حالت متنبط خودمان باید episode ما تردع کرده و ررسی نع active RL rsb behavioral Poligibis , 2 = passive RL ودر حالت كر طبق أن sample جمع في لغ البد الريست كرى لئع في وان آزارہ را علما کرنے جون ملی بع سے کے بین آیا ماکای در معط بلغ با دائی تعدادی Scriple را یک مخص کرد کخص دیگری قبل جمع کردن ایت می توان می در آن حالت نقع منفص کری مَلِا بِ عَلَى اَنه بِن اَن عَمِل اِنه مِن اَن عَلِي مَلِل بِ مِن اَورده Perard, la state de la State de C. عا من مع بات دری صورت مع بالیسی کی بالت بهذات

جی غلط ، آین الگوریع فارنج از تنظین زدن T و R متی



a - crg max (5,6) (4) (P  $\begin{array}{c} (S_1) = EL & 2 = ES \\ (S_2) = EL \\ (S_2) = EL \\ (S_3) = EL \\ (S_4) = EL /ES \\ (S_4)$ (Se)= Es 2 (SE) = Es ع کارکو نرف یکند که هر حالت منه تا بع حاله میلی خود و به هم حالات مَلِی بَعَی ندارد ر این جون تعدار State هار ا کو کردع طیعتا متداری داره را از از حاده امع وسلمان امل کامل تراست از طی اکو ہے تنا رے دارنے مل آر اکمے عاد ماری کے دا درای کھالت ماری کی تناد کاملا وا منع اے در واقع درحالہ دور استقال حالہ عای گذیرہ آیدی کو State اور ارزار نیت



$$= \frac{R_1(s,a)}{R_2(s,a)} R_{(s,a)} \times P_{(s)} P_{(a|s)} =$$

$$\leq R(s,\alpha)P(s,\alpha)(s,\alpha) = E[R(s,\alpha)]$$

$$= s - P(s), \alpha + 2 (s,\alpha)$$

$$J = \frac{1}{h} = \frac{1}{2} \frac{2}{2} (5,a) R(5,a)$$

$$\hat{J} = E_{S \in P(S), \alpha = 2, (S, \alpha)} \left[ \frac{2(S, \alpha)}{2(S, \alpha)} R_{(S, \alpha)} \right]$$

$$\lim_{N\to\infty} \left( \left| \hat{J} - \frac{1}{2} \left[ \frac{R(s, a)}{(s, a)} \right] \right| > E \right) = 0$$

$$\lim_{N\to\infty} \left( \left| \hat{J} - \frac{1}{2} \left[ \frac{R(s, a)}{(s, a)} \right] \right| > E \right) = 0$$

$$\lim_{N\to\infty} \left( \left| \hat{J} - \frac{1}{2} \left[ \frac{R(s, a)}{(s, a)} \right] \right| + \left| \frac{R(s, a)}{(s, a)} \right| + \left| \frac{R(s,$$

$$\lim_{\lambda \to \infty} \frac{1}{\sum_{i=1}^{n}} \frac{2_{i}(s_{i}a)}{2_{o}(s_{i}a)} R(s_{i}a) = 0$$

$$\frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_1(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{R_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s), \omega_2(s, \omega)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s)} \left[ \frac{R_2(s, \omega)}{R_2(s, \omega)} \right] = \frac{E}{S \sim P(s)} \left[ \frac{$$

$$\lim_{h\to 0} \frac{1}{h} \approx \frac{2}{2} (s_{20}) = \left[ \frac{2}{2} (s_{20}) \right]$$

$$= \frac{Z^{2}((5,0))}{2_{o}(5,0)} P P = \frac{Z^{2}((5,0))}{2_{o}(5,0)} P P = \frac{Z^{2}((5,0))}{2_{o}(5,0)}$$

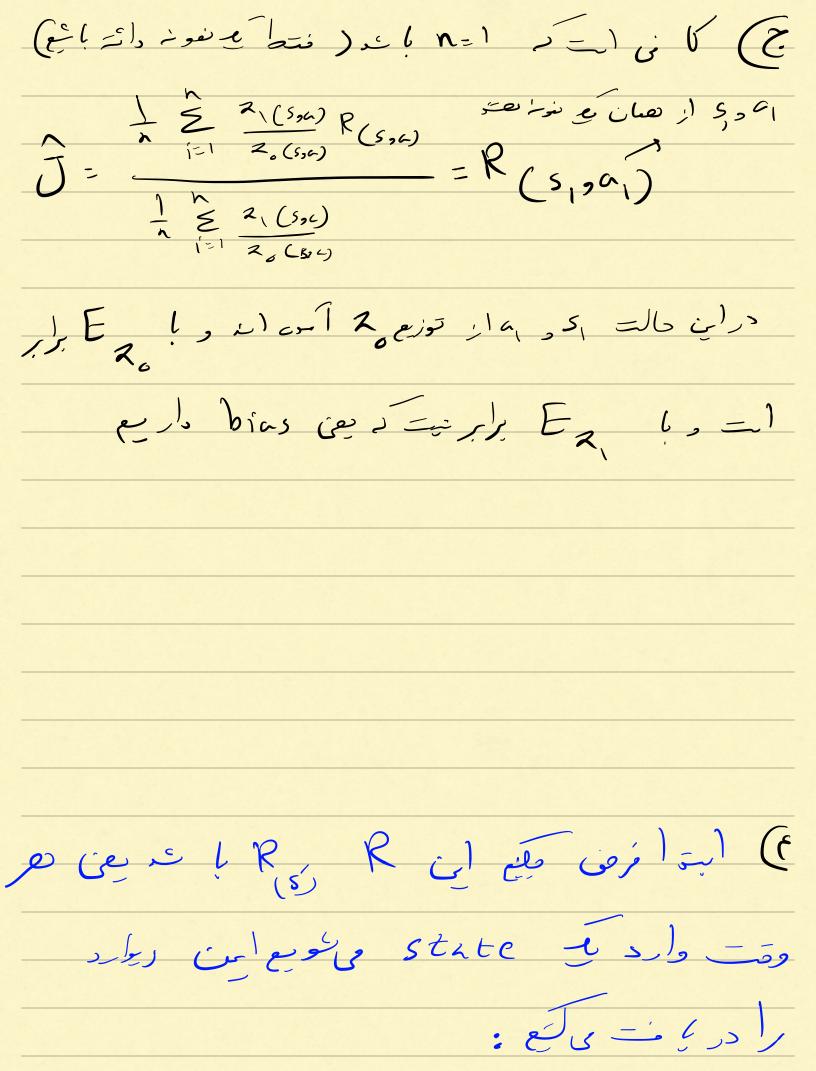
= 
$$\sum 2(590)^{2}(5) = 1 - 1000$$

$$= \sum_{s=p(s),s\in\mathbb{Z}_{1}(s,a)} \left[ \mathbb{R}(s,a) \right]$$

$$= \sum_{s=p(s),s\in\mathbb{Z}_{1}(s,a)} \left[ \mathbb{R}(s,a) \right]$$

$$= \sum_{s=p(s),s\in\mathbb{Z}_{1}(s,a)} \left[ \mathbb{R}(s,a) \right]$$

سے این مورد نیز تان داری کو برار نصان اور را می صور =



V . 2 = Peace

الن

$$\gamma_{(5)}^{2} = \sum_{s} T(s, z_{(5)}, s) \left[ R(s, z_{(5)}, s) + \gamma_{(s')}^{2} \right]$$

$$\rightarrow \frac{1}{r} \left[ r + \gamma \gamma_{(R)}^{2} \right] + \frac{1}{r} \left[ -1 + \gamma \gamma_{(D)}^{2} \right]$$

$$\gamma_{(R)}^{2} = \left[ \Gamma + \gamma \gamma_{(R)}^{2} \right] \rightarrow \left[ \gamma_{(R)}^{2} = \frac{\Gamma}{1 - \gamma} \right]$$

$$\gamma_{(D)}^{2} = \left[ -1 + \gamma \gamma_{(D)}^{2} \right] \rightarrow \left[ \gamma_{(D)}^{2} - \frac{1}{1 - \delta} \right]$$

$$\frac{1}{\sqrt{m}} = \frac{1}{\sqrt{n}} \left[ \frac{1-\delta}{1-\delta} \right] = \frac{1}{\sqrt{n}} \left[ \frac{1}{\sqrt{n}} + \frac{1}{\sqrt{n}} \right]$$

$$=\frac{1}{1-\delta}-\frac{1}{\tau}$$

$$\frac{1}{1-\delta}$$

$$\frac{2}{\tau(1-\delta)}$$

$$\gamma_{(m)} = + \Omega, \gamma_{(0)}^2 = -10, \gamma_{(R)}^2 = 70$$

· ce e is policy improvement chia il s

$$\frac{2}{i+1}(s) = \underset{A}{\operatorname{argmax}} \quad Q^{2}(s, \omega) \quad \underset{War \to W}{\operatorname{Peace}} \rightarrow P$$

$$y \quad Q^{2}(s, \omega) = \underbrace{\sum}_{s} T(s, \alpha, s) \left[ R(s, \alpha, s) + 8 \tau_{(s)}^{2} \right]$$

$$\Rightarrow \underbrace{2(s)}_{i+1} = \underset{A}{\operatorname{argmax}} \underbrace{\sum}_{s} T(s, \alpha, s) \left[ R(s, \alpha, s) + 8 \tau_{(s)}^{2} \right]$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) \left[ R(s, \alpha, s) + 8 \tau_{(s)}^{2} \right]$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) \left[ R(s, \alpha, s) + 8 \tau_{(s)}^{2} \right]$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{1}_{1e} \times G = \alpha | \Omega$$

$$= \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{\sum}_{i+1} (s, \alpha, s) + 4 \underbrace{\sum}_{i+1} (s, \alpha, s) + 8 \tau_{(s)}^{2} = 1 + \underbrace{\sum}_{i+1} (s, \alpha, s) + 4 \underbrace{\sum}_{i+1} (s, \alpha, s) +$$

$$\varphi^{2}(p,p) = -16, \quad \varphi^{2}(p,w) = \omega/\omega$$

$$\Rightarrow \begin{bmatrix} z^{i+1} - war \\ (p) \end{bmatrix} = \omega$$

$$Q(D, w) = 0 + \frac{1}{r} [-r + 0 - 0] = -1$$

sample 2:

Sample 4.

M-P	R-P	D-W
6	6	0
0	0	_1
0	0	+
0	clas	+1
6/16	0/6	+1

معی دین از state کارج و لوع این ربوارد

√ . 2 = Peace

ال

M- mountain R- Riverside D. Desert

$$\gamma_{(3)}^{2} = \sum_{s} T(s, z_{(s)}, s) \left[ R(s, z_{(s)}, s) + \gamma_{(s')}^{2} \right]$$

$$\rightarrow \frac{1}{r} \left[ 1 + \gamma \gamma_{(R)}^{3} \right] + \frac{1}{r} \left[ 1 + \gamma \gamma_{(0)}^{3} \right]$$

$$\gamma_{(R)}^{2} = \left[ \Gamma_{+} \gamma \gamma_{(R)}^{2} \right] \rightarrow \left[ \gamma_{(R)}^{2} = \Gamma_{(R)}^{2} \right]$$

$$\gamma_{(D)}^{2} = \left[ -1 + \gamma \gamma_{(D)}^{2} \right] \rightarrow \left( \gamma_{(D)}^{2} - \frac{1}{1 - \delta} \right)$$

$$\frac{1}{2} + \frac{1}{2} = \frac{1}{2} \frac{1}$$

$$=\frac{Y-Y}{Z-Y}$$

$$\frac{2}{Y-Y}$$

$$\frac{1}{Y-Y}$$

$$\gamma_{(m)} = \omega/\omega, \ \gamma_{(0)}^2 = -10 \ \gamma_{(R)}^2 = 70$$

: ce e is policy improvement die ils

2 (s) = argmax 
$$Q^{2}i(s, a)$$
 Peace  $\rightarrow P$ 

i+1 a  $Var \rightarrow V$ 

$$Q_{(N,P)} = \frac{1}{7} [1 + \frac{9}{16} \times 6] + \frac{1}{7} [1 + \frac{9}{16} \times -16] = 0/6$$

$$Q(M,W) = \frac{1}{16} \left[ 1 + \frac{9}{16} \times 40 \right] + \frac{5}{16} \left[ 1 + \frac{9}{16} \times -16 \right] + \frac{5}{16} \left[ 1 + \frac{9}{16} \times 6 \right]$$

$$\frac{1}{2^{i+1}} = war$$

