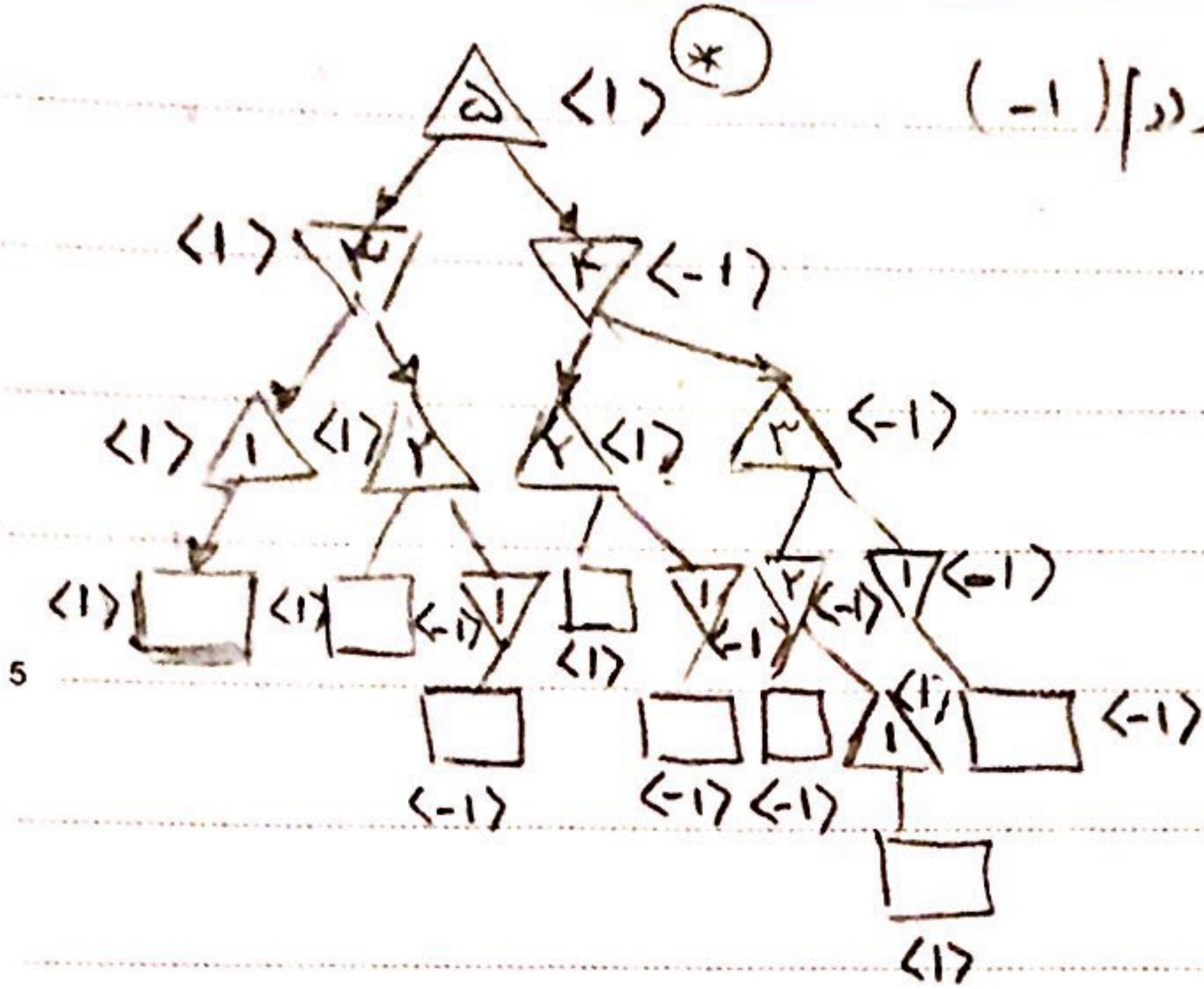
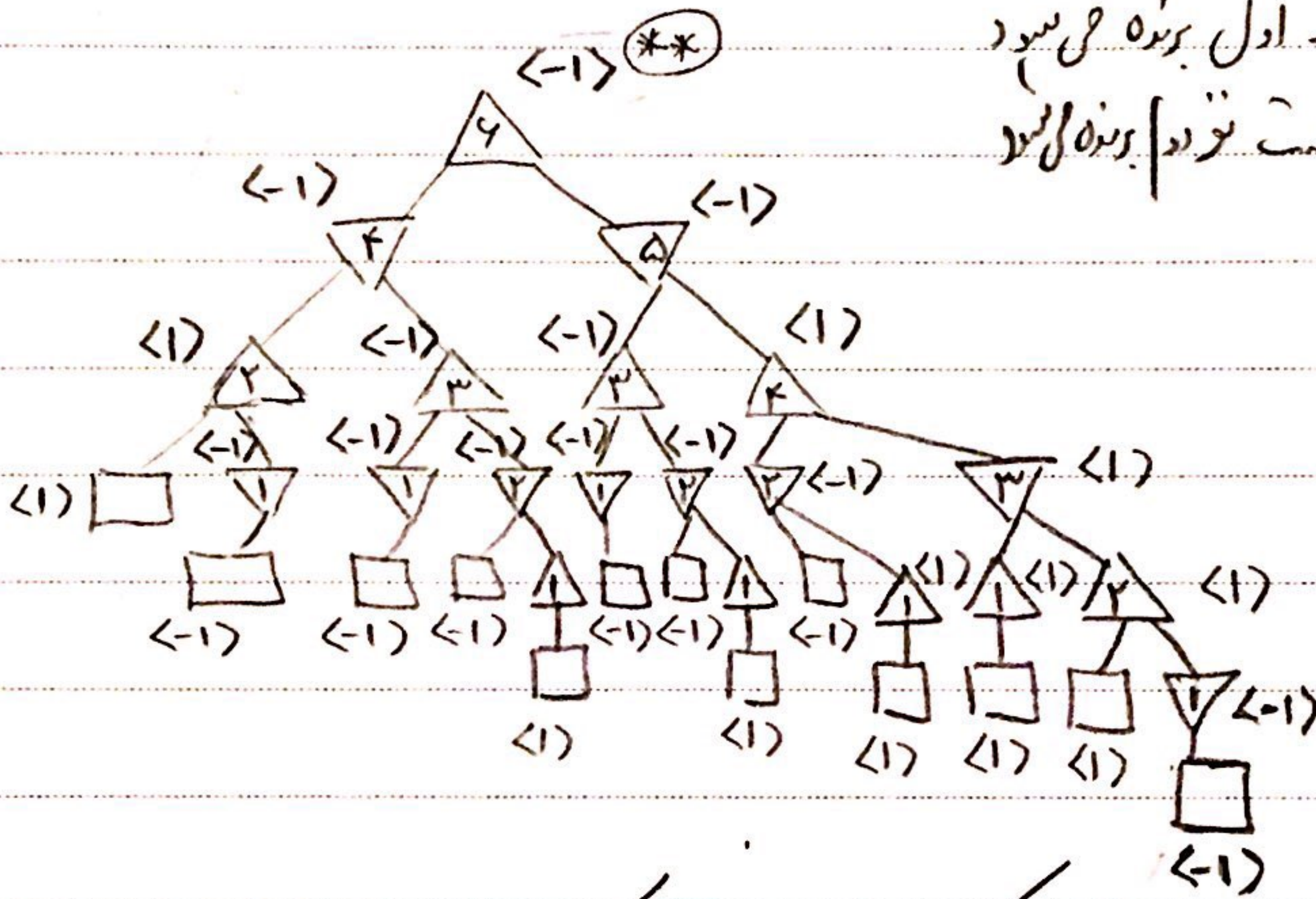


سوال ۱) آ) امتیازهای هر گره بر حسب برد توالی (۱) یا تودم (-۱)

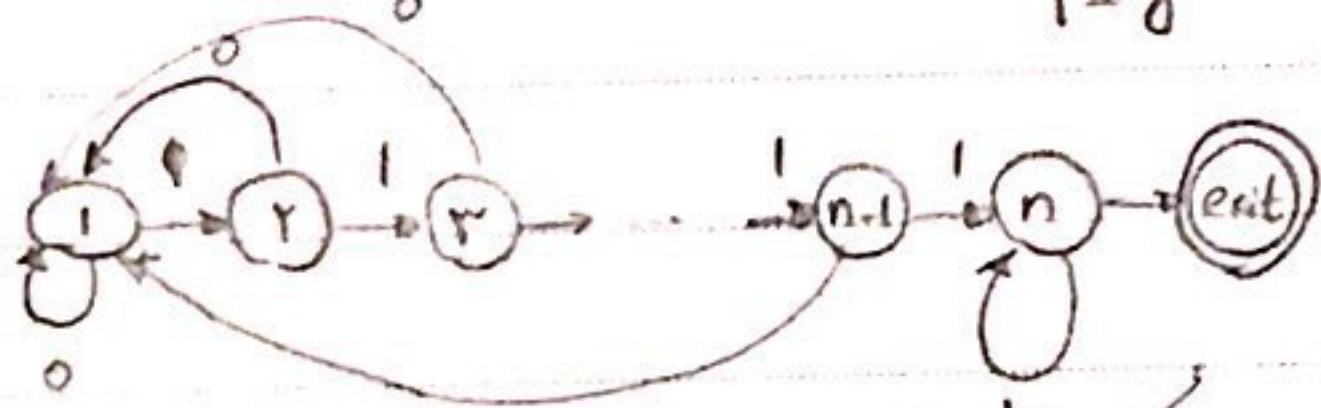
به ترتیب به صورت <۱> و <-۱> مشخص شده اند
و امتیاز داخل هر state بگیر عدد مثبت صبی می باشد



*) همان طور که مشخص است تو اول برنده می شوی
ب) همان طور که مشخص است تو دوم برنده می شوی



ج) با کسی دقت به بازی می توان متوجه شد که اعدادی وجود دارند که هر فرد در تلاش برای رسیدن به آن (پس از هوشی) هستند تا بتوانند تو مقابل را ببرند. برای مثال اگر تو اول بتوانی با معادلی که هوش عدد را به ۳ برسانی برد از قطعی است لذا در صورتی که عدد نوبت او ۴ باشد او می تواند با ۳ برنده شود. حال اگر عدد نوبت او ۶ باشد او می تواند عدد را به ۴ یا ۵ که هوش دهد که در صورت تو دوم می تواند با ۳ برنده شود. لذا استراتژی برد این گونه است که هر فرد تلاش کند تا با ۳ عدد به عددی منفی ۳ برساند و برای عدد خرابی کند. اگر عدد ابتدای آن منفی ۳ باشد تو دوم کت بازی بسته برنده می شوی و در غیر این صورت تو اول برنده می شوی



سوال ۲ الف

همانطور که از شکل بالا مشخص است در هر مرحله جازه حاصل از برکت به خانه ۱ برابر صفر است و این را ادامه در صورت

بازگشت با جازه حاصل یک پیشروی در مسیر مواجه خواهیم شد که ضریب ۸ می گیرند و این فریب به صورت توان کاهش می یابد

این در حالی است که رفتن به وضعیت n و ماندن در این وضعیت جازه ۱۰ را به ما می دهد که با توجه به جازه های

یک قبلی مشخص است هر چه سریع تر خود را به وضعیت n برسانیم و در این وضعیت بمانیم. هدف از حوصله سریع تر رسیدن به

وضعیت n این است که در هر مرحله ۸ باعث کاهش جازه می شود و به نوعی بهتر است زودتر به جازه ۱۰ برسیم.

لذا Policy مورد نظر این است که در مراحل $n-1$ به استیتم شماره بعد برویم و در حالت استیتم n در این وضعیت

بمانیم. البته بدان این Policy با توجه به توضیحات داده شده منس بر تلاش برای رسیدن به جازه ۱۰ در سریع ترین حالت است که می شود

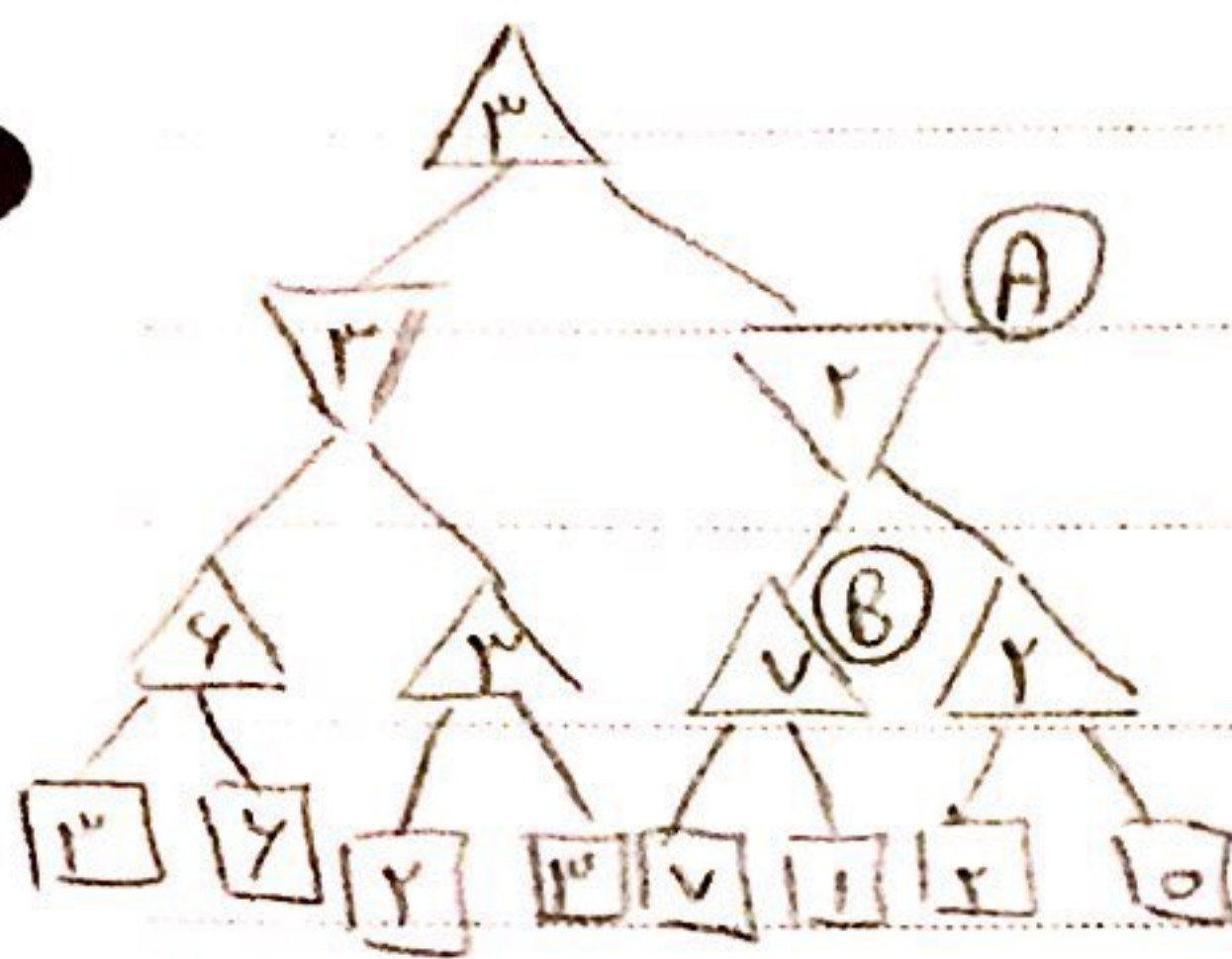
$$V^*(n) = 10 + 8V^*(n) \Rightarrow V^*(n) = \frac{10}{1-8} = 20 \quad (ب)$$

$$V^*(n-1) = 1 + 8V^*(n)$$

$$V^*(n-2) = 1 + 8V^*(n-1) = 1 + 8 + 8^2 V^*(n) \Rightarrow V^*(n-i) = \sum_{j=0}^{i-1} 8^j + 8^i V^*(n)$$

$$i \leq n-1 \Rightarrow V^*(i) = \sum_{j=0}^{n-i-1} 8^j + 8^{n-i} V^*(n) = \frac{1-8^{n-i}}{1-8} + 8^{n-i} \times 20 = 2 - \left(\frac{1}{8}\right)^{n-i-1} + \left(\frac{1}{8}\right)^{n-i} \times 20$$

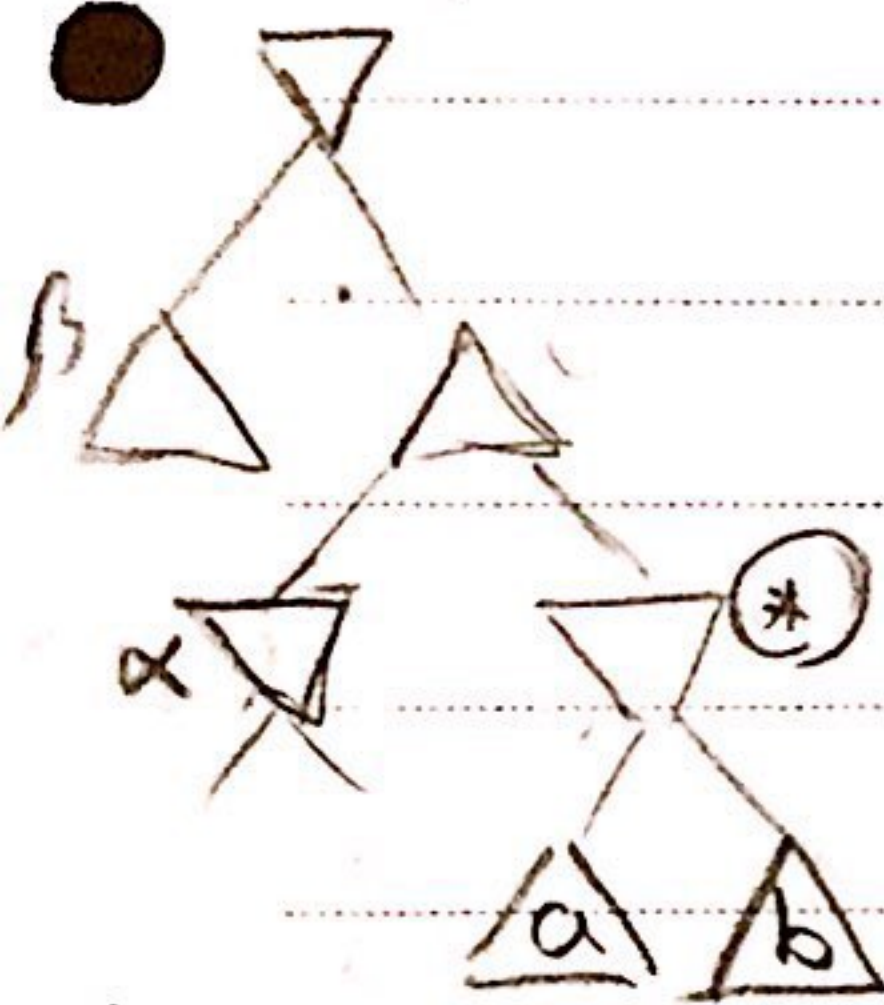
سوال ۳) آ) ابتدای وقت نمی فهمی را تحلیل می کنیم



بهترین امیدی که می تواند این برنده شود ۷ است به این گونه که از قصد
که (A) را انتخاب می کند و سپس با ریشه دادن به قصد
موجب این می شود که سبک تره (B) را انتخاب کند و لذا

این ۷ امتیاز برنده می شود که لذا میزان دلار برنده شود و با حساب هزینه فساد برابر $7 - 2 = 5$ خواهد بود.

ب) اکنون که هزینه منوره به ۵ دلار افزایش یافته است می توان فهمید که دلار بیشتر با حساب رسیدن به بهترین
امتیاز با کمک یک روش $7 - 5 = 2$ خواهد بود که از دلار این در حالت عادی کمتر است لذا بهتر است
با توجه به قیمت جدید ریشه این به صورت سالم بازی کند و ۳ دلار برنده شود.



ج) در حالت عادی برای بعضی از گره ها بر پایه حالت بدود
انتخاب می شود.

α و β $(*)$ is pruned

اما در این سوال آمدن دارد میزان با ریشه دادن ای قبل هر حالت $(*)$ به جای $\min(a, b)$ و $\max(a, b)$ با انتخاب
کرد و این pruning به اشتباه برداشت حالت $(*)$ را prune می کند که این که موجب از دست رفتن امکان رسیدن
به بهترین دلار با توجه به امکان ریشه دادن می شود

$$Q(s, a) = \sum_{s'} T(s, a, s') (R(s, a, s') + \gamma V(s'))$$

$$\Rightarrow Q(s, a) = R(s, a, s') + \gamma V(s')$$

(سوال ۴) الف)

چون داریم در این سوال $T(s, a, s') = 1$ چاره بریزیم

$$\Rightarrow Q(s, left) = 100 + \gamma V(s')$$

$$V(s') = R(s, left, s') + \gamma V(s'')$$

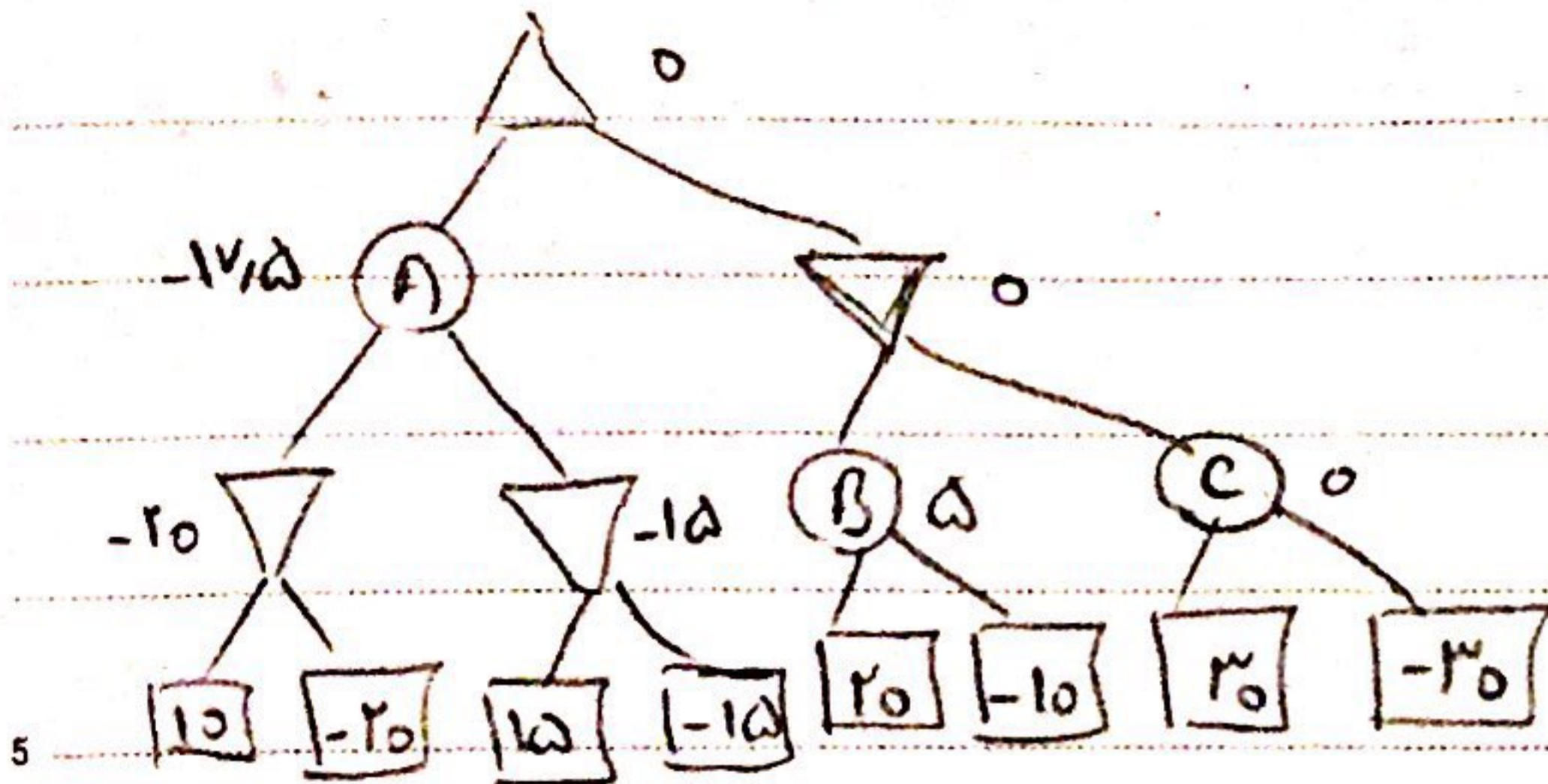
از آن جا که برای همه s policy مشخص است داریم :

$$\begin{aligned} \Rightarrow Q(s, left) &= 100 + \gamma (-1 + \gamma (-1 + \gamma (\dots))) = 100 - \gamma - \gamma^2 - \dots - \gamma^{99} \\ &= 100 - \gamma (1 + \dots + \gamma^{99}) \approx 100 - \frac{\gamma}{1-\gamma} \end{aligned}$$

$$Q(s, down) = 50 + \gamma (1 + \gamma (1 + \gamma (\dots))) = 50 + \gamma + \dots + \gamma^{99} \approx 50 + \frac{\gamma}{1-\gamma}$$

$$Q(s, down) > Q(s, left) \Rightarrow \frac{\gamma}{1-\gamma} > 25 \Rightarrow 24\gamma > 25 \Rightarrow \gamma > \frac{25}{24}$$

(—)



$$E[\text{root}] = 0$$

سوال ۵) درخت شجره‌ای را تکمیل کنید

$$E[B] = 20 \times \frac{1}{2} + (-10) \times \frac{1}{2} = 5$$

$$E[C] = 30 \times \frac{1}{2} + (-30) \times \frac{1}{2} = 0$$

$$E[A] = -20 \times \frac{1}{2} + (-15) \times \frac{1}{2} = -17.5$$

ب) به ازای هر حالت معلوم بران هر سله مقدار رتبه را می‌نویسیم

	چپ	راست
A	0	0
B	0	-10
C	5	-17.5

پس به جدول اهمیت داشتن C می‌اندازیم چرا که اگر C را انتخاب کنیم نود اول حاصل A را انتخاب می‌کند و اگر C چپ باشد نود اول همان شاقه سمت راست را انتخاب می‌کند و این را می‌توانیم در نظر بگیریم. لذا نتیجه C در انتخاب نود اول مد نظر قرار می‌دهد این در حالی است که اطلاع از نتیجه A برای نود اول در انتخاب نود اول ندارد چرا که او همیشه شاقه سمت راست را انتخاب می‌کند زیرا می‌توانیم که راست از A بیشتر است ۵ > ۱۷.۵

۲) مشخص است بین حالت‌های $E[\text{root} | A, B, C]$ بهترین حالت برای می‌توانیم رتبه

$$E[\text{root} | C = \text{چپ}] = 5$$

است لذا مشخص است که در این سله C برابر حالت چپ بهترین منفعت را برای او دارد