



پاسخ تکلیف سوم
"یادگیری تعاملی"
گروه دوم

نام استاد:

دکتر نیلی

تهیه کننده :

۸۱۰۱۹۹۲۸۹

عرفان میرزایی



پاسخ پرسش ۱:

۱. مدلسازی:

در این مسئله باید ابتدا MDP مورد نظر برای سوال را تعریف کنیم. ما در این مدلسازی حالت را در هر مرحله یک چهارتایی مرتب شامل قیمت هر یک از سهام ها به ترتیب و در آخر میزان دارایی عامل در نظر گرفتیم. که ارزش سهام شرکت ها عددی بین ۱ تا ۱۰ است (میزان ارزش سهام شرکت به دلار تقسیم بر ۵) و دارایی عامل عددی از بین ۱ تا ۲۰ (یعنی از ۵ تا ۱۰۰ دلار) هم چنین عملی که در هر مرحله می توانیم انجام دهیم را به صورت یک سه تایی مرتب در نظر گرفتیم که هر کدام از مولفه های آن می تواند صفر یا یک باشد. عملی که همه ی مولفه های آن یک باشد را خارج از مجموعه ی اعمال در نظر می گیریم.

مدل ارائه شده ویژگی مارکوف را دارد به این معنا که وقتی در یک حالت قرار داریم با فرض انجام یک عمل مشخص، حالت بعدی و پاداش دریافتی کاملاً مستقل از حالت ها و اعمال ما در زمان های گذشته است.

حال باید احتمال گذار هر یک از حالت ها به حالت دیگری را تعریف کنیم که این احتمال ها با توجه به جدول های داده شده در صورت سوال و هم چنین ارزش اولیه و ثانویه سهام که از حالت ها بدست می آید با فرض مستقل بودن تغییر قیمت شرکت ها از یکدیگر به سادگی قابل محاسبه است. هم چنین با توجه به مدل پیاده سازی شده برای گذار از یک حالت به حالت دیگر تنها یک پاداش وجود دارد که مقدار آن از تفاوت دارایی های عامل در دو حالت بدست می آید.

نکات و فرضیات :

۱. در این مسئله در نظر گرفته شده است که عامل در ابتدای هر مرحله ی زمانی سهام را خریده و در پایان آن روز سهام را در هر صورتی می فروشد و دارایی عامل به روز رسانی می شود.
۲. با توجه به صورت مسئله حالت های نهایی را زمانی در نظر گرفته ایم که اگر دارایی عامل به ۱۰۰ یا ۵ دلار برسد مسئله خاتمه می یابد و در مدل ما در این صورت به ترتیب عامل پاداشی برابر ۱۰۰+ و یا - ۱۰۰ دریافت می کند که معادل ۵۰۰+ و ۵۰۰- دلار می باشد.
۳. در این مسئله با توجه به آنکه احتمال تغییر ارزش سهام با توجه به مقدار آن می تواند متفاوت باشد (در حداکثر و حداقل ارزش) مجبوریم که این مقادیر را جز حالت های عامل در نظر بگیریم.
۴. در این مدل سازی در نظر گرفته شده است که اگر عامل عمل خرید سهامی را انتخاب کرد که بیشتر از دارایی آن است ارزش شرکت ها با توجه به احتمالات تغییر کرده ولی دارایی عامل ثابت می ماند.



۲. پیاده سازی :

پس از مدل سازی مسئله در قالب MDP باید آن را پیاده سازی کنیم. برای حل یک MDP یعنی بدست آوردن سیاست بهینه مجبور به استفاده از برنامه نویسی پویا هستیم که در این سوال از روش Policy iteration استفاده می کنیم. به این منظور از شبه کد زیر که مربوط به کتاب ساتون و بارتو است استفاده شده است.

Policy Iteration (using iterative policy evaluation) for estimating $\pi \approx \pi_*$

1. Initialization
 $V(s) \in \mathbb{R}$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$
2. Policy Evaluation
Loop:
 $\Delta \leftarrow 0$
 Loop for each $s \in \mathcal{S}$:
 $v \leftarrow V(s)$
 $V(s) \leftarrow \sum_{s',r} p(s', r | s, \pi(s)) [r + \gamma V(s')]$
 $\Delta \leftarrow \max(\Delta, |v - V(s)|)$
 until $\Delta < \theta$ (a small positive number determining the accuracy of estimation)
3. Policy Improvement
 policy-stable \leftarrow true
 For each $s \in \mathcal{S}$:
 old-action $\leftarrow \pi(s)$
 $\pi(s) \leftarrow \operatorname{argmax}_a \sum_{s',r} p(s', r | s, a) [r + \gamma V(s')]$
 If *old-action* $\neq \pi(s)$, then *policy-stable* \leftarrow false
 If *policy-stable*, then stop and return $V \approx v_*$ and $\pi \approx \pi_*$; else go to 2

شکل ۱- شبه کد Policy Iteration

که به این منظور کلاس محیط StockMarketEnvironment را پیاده سازی می کنیم. این کلاس از EnvironmentBase مشتق شده است و ویژگی های آن را ارث بری کرده است. مجموعه حالت ها و اعمال به صورت گفته شده در بخش قبلی در نظر گرفته شده است که بدین منظور از MultiDiscrete کتابخانه ی gym استفاده کردیم. مقدار نرخ تخفیف و هم چنین احتمال گذارها را به عنوان ورودی دریافت می کنیم.



شبه کد پیاده سازی شده را در متد `Policy Iteration` پیاده سازی کرده ایم که این متد یک تابع ارزش اولیه و یک سیاست قطعی را به عنوان ورودی دریافت می کند و سیاست و تابع ارزش بهینه را به عنوان خروجی باز می گرداند. هم چنین مقدار تتا را نیز به عنوان ورودی دریافت می کند که دقت ارزیابی ارزش ها را مشخص می کند.

با استفاده از متد `get_all_states, get_available_actions` تمام حالت ها و عمل های ممکن را در خروجی باز می گردانیم. برای پیاده سازی این شبه کد از یک حلقه `While` بزرگ استفاده کردیم که تا وقتی سیاست پایدار نشده یعنی به ازای انجام مرحله ی بهبود سیاست، سیاست مربوط برخی از حالت ها تغییر می کند این کار را ادامه می دهیم. بدین صورت که در هر مرتبه ابتدا ارزش ها را به روزرسانی کرده تا زمانی که بیشترین اختلاف میان دو ارزش یک حالت به کمتر از تتا برسد. برای به روز رسانی ارزش یک حالت می بایست که احتمال گذار از یک حالت به حالت بعدی و هم چنین پاداش دریافتی در آن حالت را داشته باشیم. که با توجه به آنکه مدل ما دارای ۲۰۰۰۰ حالت است بدین صورت عمل کرده ایم که تنها این مقادیر را برای حالت های همسایه یک حالت محاسبه کردیم و این کار را در متد `get_dynamics` پیاده سازی کرده ایم. بدین صورت که در هر حالت-عمل حالت های همسایه و احتمال رفتن به این حالت ها و هم چنین پاداش گذار به این حالت را به عنوان خروجی باز می گردانیم. که برای تولید حالت های مجاور، از شروط استفاده می کنیم.

هم چنین پس از بدست آوردن این احتمالات و پاداش ها باید مقدار ارزش هر حالت را به روز رسانی کنیم که این کار را به کمک متد `update_values` انجام می دهیم که به ازای همه ی حالت های مجاور مقدار احتمال را در مجموع پاداش و ارزش کاهش یافته حالت مجاور ضرب می کند و مقادیر آن ها را با یکدیگر جمع می کند.

پس از پیاده سازی متد `Policy Iteration` در نوت بوک `Q1.ipynb` که مربوط به این سوال است ابتدا کتابخانه های لازم را وارد می کنیم. سپس احتمال های تغییر ارزش سهام شرکت ها و هم چنین عمل ها را به عنوان ورودی های لازم برای پاس دادن به محیط تعریف می کنیم. هم چنین یک سیاست قطعی که در هر حالت به تصادف یک عمل را انتخاب می کند را به همراه تابع ارزش صفر برای هر حالت به عنوان ورودی به الگوریتم `Policy Iteration` می دهیم.

۳. نتایج :

برای با خبر شدن از میزان پیشرفت الگوریتم در هر بار اجرای حلقه ی خارجی الگوریتم به ازای هر بار که الگوریتم `Value Evaluation` اجرا می گردد میزان دلتا که با تتا مقایسه می شود را چاپ می کنیم که همان طور مشاهده می شود در هر بار تلاش این مقدار کاهش می یابد و هم چنین برای هر بار اجرای حلقه خارجی نیز این مقادیر



کوچکتر از مقادیر گذشته است که این (کوچک تر شدن تغییرات ارزش ها) می تواند نشانه ای از نزدیک به همگرا شدن الگوریتم باشد.

هم چنین در مرحله Policy Improvement با استفاده از یک شمارنده تعداد حالت هایی که سیاست آن عوض شده اند را حساب کرده و چاپ می کنیم. که همان طور که انتظار داشتیم این مقدار برای اولین مرحله بدلیل تصادفی بودن مقدار بسیار بالایی دارد (۱۸۰۰۰) و رفته رفته و در طی پنج گام این تغییرات به صفر می رسد که این امر نشان دهنده ی این است که با توجه به ارزش های حالات که با دقت تتا بدست آمده است توانسته ایم که سیاست بهینه را بدست بیاوریم.

تتا در تمامی مراحل به مقدار ۰,۰۲ محاسبه شده است که این مقدار با توجه به اینکه ارزش پاداش ها مقداری بین ۱۰۰+ و ۱+ و ۰ و ۱- و ۱۰۰- است مقدار کوچکی به محسوب می شود.

با توجه به اینکه در این مسئله احتمال گذار بین حالت های مختلف را داشتیم می توانیم امید ریاضی هر عمل را محاسبه کنیم با این فرض که در هر حالت ما قادر به انجام تمامی اعمال هستیم. که می دانیم این فرض با توجه به اینکه ارزش شرکت ها متغیر با زمان است و همچنین دارایی عامل نیز تغییر می کند در همه ی شرایط برقرار نیست.

$$E[\text{Reward}(0,0,0)] = 0$$

$$E[\text{Reward}(1,0,0)] = 0.4 * (+5) + 0.3 * (0) + 0.3 * (-5) = +0.5$$

$$E[\text{Reward}(0,1,0)] = 0.1 * (+5) + 0.8 * (0) + 0.1 * (-5) = 0$$

$$E[\text{Reward}(0,0,1)] = 0.2 * (+5) + 0.1 * (0) + 0.7 * (-5) = -2.5$$

$$E[\text{Reward}(1,1,0)] = 0.04 * (+10) + 0.35 * (+5) + 0.31 * (0) + 0.27 * (-5) + 0.03 * (-10) = +0.5$$

$$E[\text{Reward}(1,0,1)] = 0.08 * (+10) + 0.1 * (+5) + 0.37 * (0) + 0.24 * (-5) + 0.21 * (-10) = -2$$

$$E[\text{Reward}(0,1,1)] = 0.02 * (+10) + 0.17 * (+5) + 0.16 * (0) + 0.57 * (-5) + 0.08 * (-10) = -2.6$$

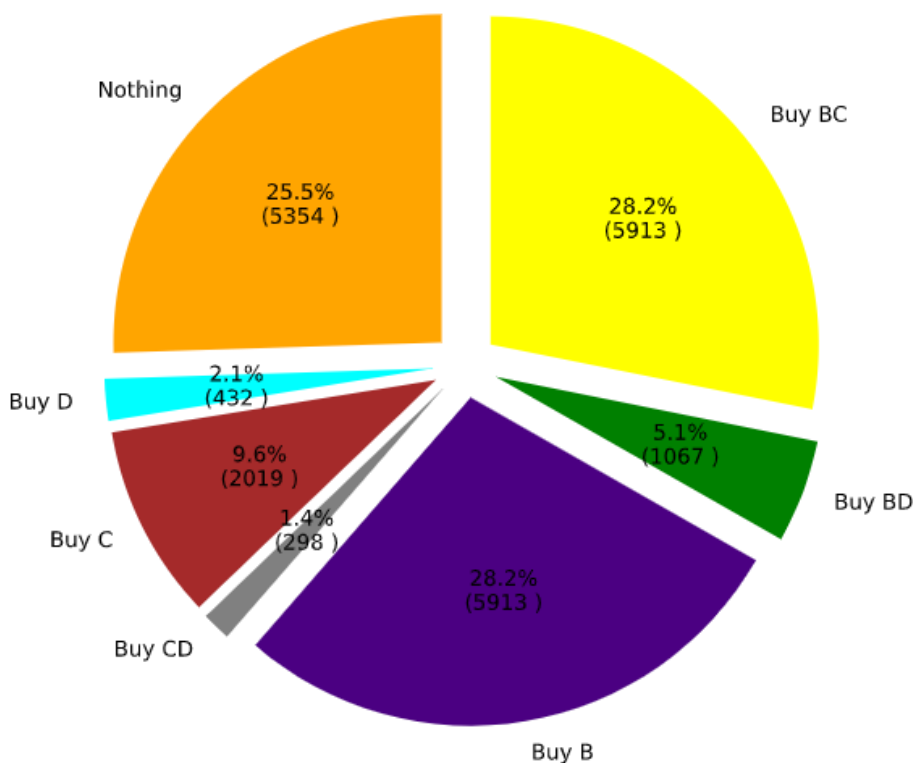
مولفه ی اول برای خرید شرکت B است و مولفه ی برای خرید سهام شرکت C و هم چنین D می باشد.



که با توجه به این مقادیر در هر مرحله در صورتی که بتواند تمام عملیات ها را انجام دهد می تواند بین حالت خرید سهام شرکت B و یا خرید سهام شرکت C , B انتخاب کند.

حال برای ارزیابی سیاست بهینه ی بدست آمده میزان انتخاب عمل های متفاوت را در یک نمودار دایره ای چاپ می کنیم.

Percentage of actions with optimal action gamma = 0.9



که همان طور که انتظار میرفت عمل خرید سهام B به همراه خرید سهام C , B به یک اندازه و در بالاترین میزان قرار دارند. که جمع آن ها تقریباً بیش از نیمی از حالت ها را به خود اختصاص می دهند که با یک محاسبه ی سر انگشتی می توان در نظر گرفت وقتی که دارایی ما ۱۱ تا ۲۰ باشد که نیمی از حالت ها را دربر می گیرد سهام شرکت ها هر چه باشد ما می توانیم تمام اعمال را انجام بدهیم و برای بقیه دارایی های از بین ۰ تا ۱۰ نیز عامل در برخی از شرایط می تواند عمل بهینه را انجام بدهد که به همین دلیل جمع آن ها بیشتر از ۵۰ درصد است.

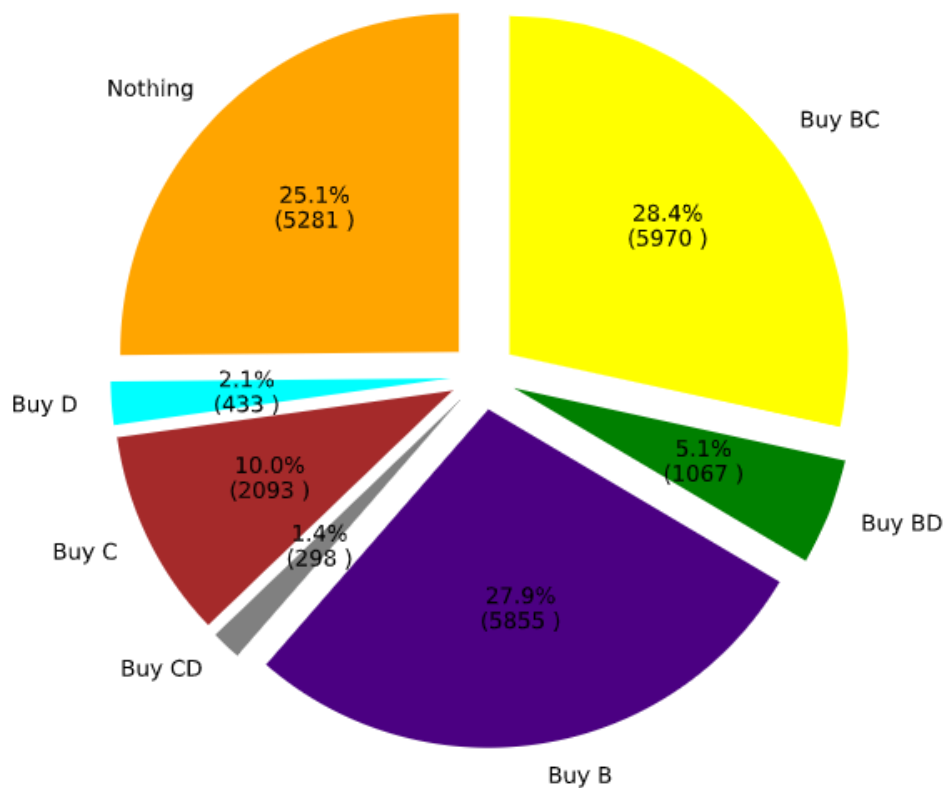


در جایگاه بعدی عمل سهامی نخریدن قرار دارد که مربوط به حالت هایی می شود که دارایی های فرد کمتر از قیمت سهام هر یک از شرکت هاست و فرد نمی تواند سهامی خریداری کند. هر چند که در این شرایط عامل می تواند هر عملی را انجام بدهد زیرا نتیجه همواره یکسان خواهد بود.

هم چنین دیده می شود که ترتیب سایر اعمال نیز دقیقاً به همان ترتیبی است که امید ریاضی آن ها را حساب کردیم و مربوط به موقعیت هایی است که به دلیل کم بودن دارایی از ارزش سهام شرکت ها می توانسته ایم هر عملی را انجام بدهیم.

۴. مقایسه ی نتایج :

Percentage of actions with optimal action gamma = 0.7

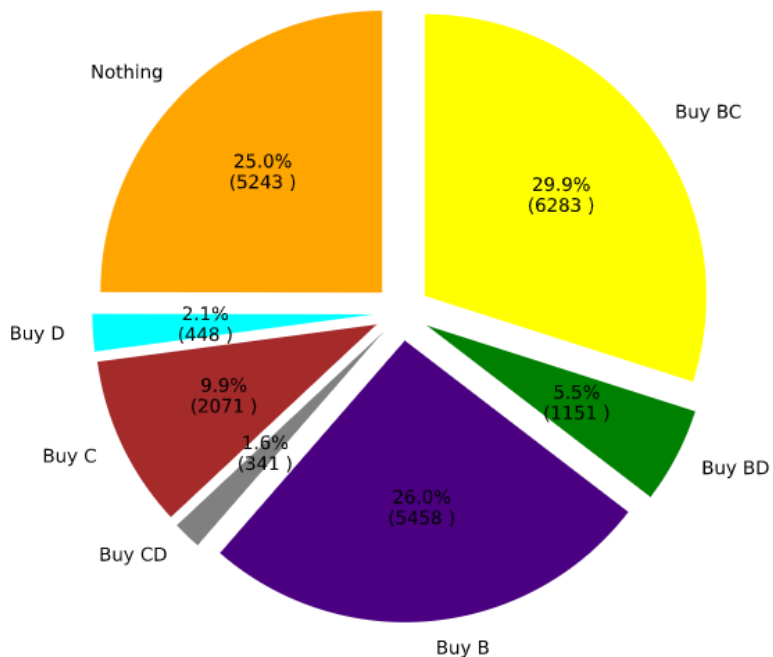




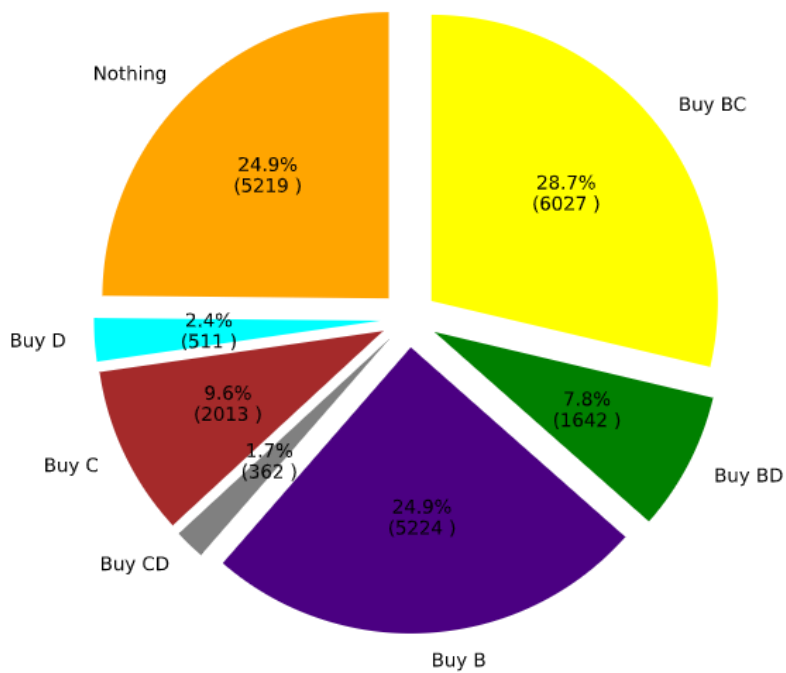
تکلیف سوم - یادگیری تعاملی
عرفان میرزایی
۸۱۰۱۹۹۲۸۹



Percentage of actions with optimal action gamma = 0.5



Percentage of actions with optimal action gamma = 0.3

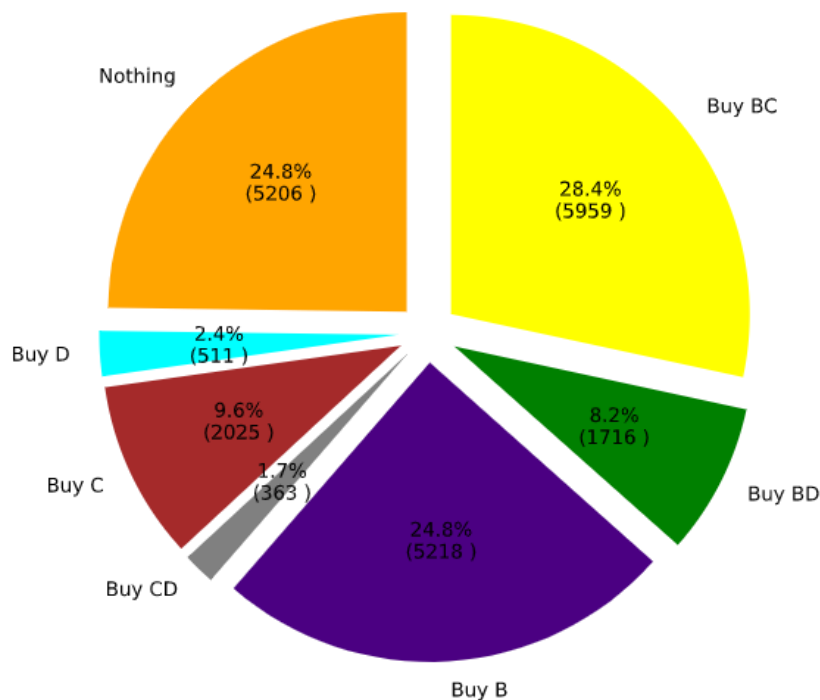




تکلیف سوم - یادگیری تعاملی
عرفان میرزایی
۸۱۰۱۹۹۲۸۹



Percentage of actions with optimal action gamma = 0.1



Gamma	۰,۹	۰,۷	۰,۵	۰,۳	۰,۱
B	۲۸.۲	۲۷.۹	۲۶	۲۴.۹	۲۴.۸
B, C	۲۸.۲	۲۸.۴	۲۹.۹	۲۸.۷	۲۸.۴
No	۲۵.۲	۲۵.۱	۲۵	۲۴.۹	۲۴.۸
C	۹.۶	۱۰	۹.۹	۹.۶	۹.۶
B, D	۵.۴	۵.۱	۵.۵	۷.۸	۸.۲
D	۲.۱	۲.۱	۲.۱	۲.۴	۲.۴
C, D	۱.۴	۱.۴	۱.۶	۱.۷	۱.۷

با توجه به اعداد ثبت شده در جدول و هم چنین نمودارهای دایره‌ای می‌توان مشاهده کرد که عملکرد عامل‌ها در تمامی مراحل نزدیک به یکدیگر است اما با کاهش میزان نرخ بهره و کوتاه‌بین شدن عامل، مقداری از انجام عمل بهینه دست می‌کشد که این اتفاق به دلیل کم شدن اثر پاداش بزرگی است که با رسیدن به دارایی ۱۰۰ می‌رسیم صورت می‌گیرد.



پاسخ پرسش ۲ :

۱. فرضیات :

۱. در این مدل سازی فرض شده است که شرکت توانایی خرید هر مقداری از محصولات اولیه را در اول ماه دارد و محدودیتی وجود ندارد. با فرض اینکه در محیط شبیه سازی برای یادگیری این مسئله مشکلی ندارد و پس از اتمام یادگیری و در دنیای واقعی عامل یاد گرفته است که اگر مادهی اولیه ای را بخرد ولی تاریخ انقضای آن تمام شود ضرر می کند.
۲. فرض کرده ایم که سهولت در تهیهی ماده اولیه و یا تاثیر تحریم بودن یا نبودن یک مادهی اولیه بر روی قیمت آن منعکس می شود و در صورت نایاب بودن مادهی اولیه ای می توان قیمت آن را به صورت حدی بی نهایت در نظر گرفت.
۳. فرض کرده ایم که مواد اولیهی لازم برای هر محصول ثابت است و با زمان تغییر نمی کند و با توجه به پاداشی که عامل در هر مرحله دریافت می کند خود می تواند به درستی یاد بگیرد که برای هر محصول چه میزان و چه مادهی اولیهی لازم است.
۴. در این مدل سازی فرض کرده ایم که سلیقهی مشتریان در ماه های مختلف سال متفاوت است اما به طور کلی، در سال های متفاوت یکسان است.
۵. با توجه به فرض گفته شده در صورت سوال که خریدار می تواند برای خرید محصولات به سایت شرکت مراجعه کند، فرض کرده ایم که خریدار پس از ورود به سایت محصولات مورد انتظار خود را وارد می کند و در صورت آنکه آن محصولات در کارخانه موجود باشد به فروش میرسد و در همان لحظه به مشتری تحویل می شود و پول آن بلافاصله دریافت می شود و همچنین در صورتی که آن جنس موجود نباشد به مشتری اطلاع داده می شود و عامل نیز پاداش منفی ای دریافت می کند.
۶. فرض کرده ایم که عامل برای یک سال برنامه ریزی شده است و در انتهای سال شروع مجدد می شود. در انتهای هر سال انبار گردانی مواد اولیه اتفاق می افتد و مواد اولیه کارخانه را صفر می کنیم و هم چنین در پایان هر ماه انبار گردانی محصولات را انجام می دهیم و تمام محصولات را صفر می کنیم.
۷. در صورتی که عامل دستور تولید تعداد محصولی را بدهد که مواد اولیهی آن در انبار نباشد، تا جایی که می توان از مواد اولیه استفاده می کنیم تا حدی که نشود یک محصول دیگر را آماده کرد.
۸. در صورتی که عامل دستور خرید محصولات غذایی را بدهد اما در روز اول ماه نباشیم هیچ تغییری در مواد اولیه رخ نمی دهد.



۹. فرض می‌کنیم در هر گام زمانی می‌توانیم تنها یک فروش از سایت انجام بدهیم.

۲. مجموعه حالت‌ها:

برای مدلسازی این مسئله با توجه به اینکه باید علاوه بر خرید مواد اولیه که در اول هر ماه اتفاق می‌افتد، تبدیل مواد اولیه به محصولات را نیز مدیریت کنیم و هم چنین فروش محصولات بر روی حالت عامل نیز تاثیر گذارند پس در نهایت گام زمانی را برابر یک دقیقه در نظر می‌گیریم.

یک ماتریس ۵۰۰ سطری در نظر می‌گیریم که در ستون اول مقدار کالای اولیه موجود در انبار است (که برای هر این اعداد مقیاسی در نظر گرفته شده است. برای مثال به گرم گزارش می‌کنیم) و در ستون دوم مدت زمان باقی مانده تا منقضی شدن آن (چند دقیقه) را نشان می‌دهد و ستون سوم آن قیمت آن ماده‌ی اولیه را بیان می‌کند. ستون اول و دوم این ماتریس برای هر سطر یک لیست شامل ۱۲ عدد است که میزان آن ماده اولیه و تاریخ انقضای هر ماده اولیه در آن روز را با توجه به ماه خریداری شده‌ی آن نشان می‌دهد.

هم چنین یک ماتریس دیگری با ۱۰۰ سطر در نظر می‌گیریم که در ستون اول تعداد محصول موجود در انبار کارخانه و در ستون دوم آن مدت زمان باقی مانده تا منقضی شدن آن (چند دقیقه) را نشان می‌دهد و ستون سوم آن قیمت آن محصول را بیان می‌کند.

ستون اول و دوم این ماتریس برای هر سطر یک لیست شامل ۱۲ عدد است که تعداد آن محصول و تاریخ انقضای هر محصول را با توجه به دقیقه‌ی تولید آن نشان می‌دهد.

حالت محیط ما در این مسئله از متشکل از ماتریس‌هایی است که در بالا ذکر شده به اضافه‌ی یک عدد ثابت که نشان‌دهنده‌ی مقدار دقیقه گذشته از ابتدای سال است.

۵. مجموعه اعمال:

این عامل در هر حالتی که قرار دارد (یعنی در هر دقیقه) می‌تواند تعداد تولید از هر محصول را مشخص کند که یک بردار ۱۰۰ تایی است و هم چنین در ادامه با یک بردار ۵۰۰ تایی میزان خرید هر ماده‌ی اولیه را (به گرم) مشخص کند. در واقع هر عمل برداری شامل ۶۰۰ مولفه‌ای است.



۶. پاداش:

در مدل سازی مسئله ها به روش یادگیری تقویتی ما با استفاده از سیگنال پاداش برای عامل مشخص می کنیم که هدف "چیست" اما در مورد "چگونگی" رسیدن به هدف اطلاعی به عامل نمی دهیم. در این مثال نیز مانند سایر مسائل یادگیری تقویتی پاداش به صورت یک اسکالر متعلق به اعداد حقیقی در نظر گرفته شده که میزان سودی است که در ماه گذشته کسب کرده ایم. هدف عامل یادگیری تقویتی این است که امید ریاضی میزان پاداش های دریافتی را در طول زمان بیشینه کند. میزان سود که در واقع به عنوان پاداش به عامل داده می شود با انجام محاسبات پیچیده ای که توسط محیط (خارج از عامل) مشخص می شود؛ بدین صورت که مقدار فروش محصولات شرکت را در دقیقه اخیر حساب کرده و سپس هزینه ی تولید این محصولات شامل مواد اولیه خریداری شده را با علامت منفی در نظر می گیریم و با آن جمع می کنیم. هم چنین میزان سودی که می توانستیم از فروش محصولاتی که مشتریان درخواست کرده بودند و ما به دلیل وجود نداشتن مواد اولیه قادر به تهیه ی آن نبودیم را نیز با علامت منفی در حاصل جمع در نظر بگیریم. هم چنین قیمت مواد اولیه و محصولاتی که تاریخ انقضای آن ها در دقیقه گذشته به پایان رسیده ولی مصرف نشده اند را با علامت منفی در جمع قرار می دهیم.

۷. نحوه ی انتقال از یک حالت به حالت دیگر:

در پایان هر دقیقه باید حالت را به روز رسانی کنیم. حالت ما شامل دو ماتریس و یک عدد بود. که آن عدد به سادگی به ازای هر دقیقه یک بار افزایش می یابد.

ماتریس مربوط به مواد اولیه در هر روز به این صورت تغییر می یابد که در ستون اول که مربوط به مقدار ماده ی اولیه موجود در انبار است، به اندازه ی مقدار مصرف شده در دقیقه گذشته به روز رسانی می شود و ستون دوم نیز هر عدد موجود در لیست به اندازه ی یک دقیقه کاهش می یابند. و قیمت ماده ی اولیه نیز ممکن است با توجه به تغییرات بازار در دقیقه گذشته تغییر کرده یا ثابت بماند.

۸. چرا MDP است ؟:

یک فرآیند تصمیم گیری مارکوف شامل یک چهارتایی مرتبط است که شامل مجموعه حالات و مجموعه عمل ها و یک تابع احتمال است که در یک حالت خاص با انجام یک عمل مشخص با چه احتمالی به حالت بعدی می رویم و هم چنین یک تابع پاداش که در یک حالت خاص با انجام یک عمل مشخص به چه میزان پاداش دریافت می کنیم.



خاصیت مارکوف بیان می‌کند که وقتی در یک حالت قرار داریم با فرض انجام یک عمل مشخص، حالت بعدی و پاداش دریافتی کاملاً مستقل از حالت‌ها و اعمال ما در زمان‌های گذشته باشد. که با توجه به توضیحاتی که در بخش قبل ارائه شد این ویژگی به طور کامل برقرار است.

هم‌چنین مجموعه حالت‌ها و اعمال می‌توانند محدود یا نامحدود باشد که با توجه به عدم فرض مسئله می‌توانیم هر کدام را در نظر بگیریم. اگر بخواهیم از مجموعه‌ی محدود استفاده کنیم باید برای قیمت‌ها سقف و گام تعیین کنیم هم‌چنین برای خرید اولیه و تولید محصول که این فرض‌ها محدودیت‌زا هستند پس از فرض نامحدود بودن مجموعه حالت و عمل استفاده می‌کنیم.

با فرض در اختیار داشتن توزیع افزایش قیمت‌ها و هم‌چنین داشتن توزیع سفارشات مشتریان و هم‌چنین با دانستن تاریخ انقضای مواد اولیه و محصولات کارخانه می‌توانیم تابع احتمال مربوط به گذار از یک حالت به حالت دیگر و پاداش دریافتی را نیز محاسبه کرد. یا اینکه می‌توانیم دینامیک محیط را ناشناخته در نظر بگیریم. که در این صورت عامل باید آن را پیدا کند.