

## کاربرد یادگیری تقویتی در یک مدل سازی عامل محور\* برای بازار عمده فروشی برق ایران

محمد رضا اصغری اسکوئی<sup>1</sup>

فرهاد فلاحی<sup>2</sup>

میثم دوستی زاده<sup>3</sup>

سعید مشیری<sup>4</sup>

تاریخ پذیرش: 1397/06/20

تاریخ دریافت: 1397/02/26

### چکیده:

مطالعات اخیر بازارهای عمده فروشی برق عموماً براساس مدل های چندعاملی است، که در آنها تعادل بازار برپایه رقابت و تعامل عوامل متعدد با یک دیگر به دست می آید. از ویژگی های اصلی این نوع مدل ها، امکان یادگیری عوامل از نتایج رفتار خود و سایرین در یک محیط رقابتی است. در بازار عمده فروشی برق، هر عامل یک واحد تولیدکننده برق است که به صورت مستقل و هوشمند با سایر عامل ها برای عرضه برق با قیمت های پیشنهادی رقابت می کند. فرآیند قیمت گذاری را می توان یک بازی ایستا فرض نمود که هرروز تکرار می شود. در این بازی هر عامل قیمت پیشنهادی خود را مستقلاً اعلام نموده و بهره بردار با توجه به تقاضای بار مصرفی و محدودیت ها، بهترین پیشنهادها را انتخاب می نماید. عامل به صورت عقلانی عمل نموده و با انتخاب استراتژی مناسب، به دنبال

\* در تهیه این مقاله از منابع و امکانات پژوهشگاه نیرو استفاده شده که به این وسیله از مسئولین و کارشناسان محترم آن مجموعه به ویژه جناب آقای دکتر کیومرث حیدری قدردانی می شود.

1. استادیار دانشکده علوم ریاضی و رایانه، دانشگاه علامه طباطبائی (نویسنده مسئول)

oskoei@atu.ac.ir

2. دانشجوی دکتری مهندسی برق قدرت، دانشگاه شاهد و پژوهشگر پژوهشگاه نیرو

ffallah@nri.ac.ir

3. استادیار دانشکده فنی و مهندسی دانشگاه لرستان

m.doostizadeh@gmail.com

4. دانشیار اقتصاد، دانشگاه ساسکاچوان، کانادا

moshiri.s@usask.ca

بیشینه نمودن سود بلندمدت خود است. در این راستا، عامل از قدرت یادگیری و بهبود استراتژی قیمت‌گذاری، که نقش بسیار تعیین‌کننده در موفقیت عامل دارد، استفاده می‌کند. یادگیری تقویتی یک روش کلاسیک است که در مدل‌های چندعاملی امکان یادگیری مبتنی بر سعی و خطا را فراهم می‌نماید. هدف این مقاله کاربرد و مطالعه روش‌های یادگیری تقویتی در مدل چندعاملی بازار برق ایران و مقایسه آن‌ها با دو استراتژی تصادفی و حریصانه است. در این مطالعه، میزان سود واحدها و زمان رسیدن به حالت تعادل به عنوان ملاک ارزیابی در نظر گرفته شده است. نتایج شبیه‌سازی نشان می‌دهد، استراتژی یادگیرنده سود عامل‌ها را به طور معناداری افزایش می‌دهد و سرعت همگرایی به حالت تعادل را بیشتر می‌کند.

### طبقه‌بندی JEL: C6, C7, Q41

کلیدواژه‌ها: مدل‌سازی عامل‌محور، بازار برق، یادگیری تقویتی، نظریه بازی‌ها، ایران

## 1. مقدمه

برق به عنوان یکی از منابع انرژی نقش به‌سزایی در رشد و توسعه اقتصادی اجتماعی کشورها داشته و میزان دسترسی به برق یکی از معیارهای توسعه یافتگی کشور به حساب می‌آید (تومن و جملکوا<sup>1</sup>، 2003). ساختار سنتی بازارهای برق با توجه به شرایط تولید آن‌ها (به عنوان نمونه برق آبی و یا برق اتمی) در بسیاری از موارد انحصاری بوده و دولت‌ها با توجه به ملاحظات زیست‌محیطی و آثار اقتصادی و اجتماعی برق به طور مستقیم یا غیرمستقیم در آن‌ها دخالت می‌کنند. در سال‌های اخیر با پیشرفت فناوری‌های جدید، بازارهای برق در بسیاری از کشورهای جهان شاهد تجدید ساختار و تحولات اساسی بوده و از یک ساختار متمرکز و یکپارچه به یک بازار آزاد و رقابتی به ویژه در بخش‌های تولید و توزیع تبدیل شده‌اند. هدف اصلی این تحول، افزایش کارآیی بازار برق و در نتیجه کاهش قیمت‌ها و افزایش رفاه مصرف‌کننده است. در عین حال، بازار عمده‌فروشی برق یک بازار نسبتاً پیچیده است و خصوصیات مانده ضرورت تعادل لحظه‌ای تولید و مصرف، محدودیت‌های تولید و انتقال، قیود فیزیکی پخش‌بار، هزینه بالای ذخیره انرژی و حساسیت پایین تقاضا نسبت به تغییرات قیمت بر پیچیدگی آن افزوده و چالش‌های خاصی در تنظیم

---

1. Toman and Jemelkoma (2003)

کاربرد یادگیری تقویتی در یک مدل‌سازی... 3

بازار ایجاد می‌کند. علاوه بر این‌ها، مسائلی مانند حمایت از محیط زیست، هزینه بالای سرمایه‌گذاری اولیه و استفاده بهینه از منابع نقش مستقیم بر عملکرد بازیگران این بازار دارد. متأسفانه در شرایط رقابتی جدید، مدل‌های سنتی و کلاسیک که عوامل اقتصادی (تولیدکنندگان یا مصرف‌کنندگان) را همگون و فاقد تعامل فرض می‌کنند کارآیی لازم را نداشته و نتایج قابل قبولی ارائه نمی‌دهند. بازارهای جدید برق شامل تعداد زیادی از تولیدکنندگان ناهمگون با توجه به نوع فناوری و درجه ریسک پذیریشان هستند که در رقابت با یکدیگر وارد چرخه تولید می‌شوند. مدل‌سازی عامل‌محور<sup>1</sup> (ABM) رفتار واحدهای تولیدی ناهمگون را تبیین و اجازه تعامل و یادگیری در یک محیط پویا را فراهم می‌سازد، اما به علت وسعت مدل و پیچیدگی‌های زیاد آن نمی‌توان از یک راه حل تحلیلی برای به دست آوردن متغیرها در شرایط تعادلی استفاده کرد. با توجه به پیشرفت فناوری محاسبات، نتایج مدل‌های عامل‌محور را معمولاً با روش‌های شبیه‌سازی در چارچوب سناریوهای مختلف می‌توان ارزیابی کرد.

به طور مشخص‌تر، بازار عمده‌فروشی برق از تعداد زیادی نیروگاه‌های بزرگ و کوچک که در قالب شرکت‌های تولیدی فعالیت دارند در کنار شرکت‌های انتقال و توزیع تشکیل شده‌است و نهاد قانون‌گذار و نهاد مستقل بهره‌بردار سیستم<sup>2</sup> (ISO)، مسئولیت وضع قوانین و کنترل و هدایت بازار را برعهده دارند. واحدهای تولیدکننده روزانه برنامه تولید پیشنهادی خود را برای هر ساعت به نهاد بهره‌بردار ارائه می‌دهند. این نهاد با توجه به پیش‌بینی تقاضای مصرف بیست و چهار ساعت آتی و با اجرای مکانیزم حراج، واحدهای برنده و قیمت بازار را مشخص می‌نماید و سپس با اجرای مکانیزم تسویه، فروش برق و تسویه حساب انجام می‌گیرد. با توجه به ضرورت تعادل مصرف و تولید در کل کشور و همچنین توزیع جغرافیایی واحدهای تولیدی و مصرف‌کننده، در مکانیزم حراج، علاوه بر میزان مصرف، قیود اقتصادی و قیود فنی مربوط به تولید، انتقال و توزیع نیز

---

1. Agent Based Model (ABM)

2. Independent System Operator (ISO)

لحاظ می‌شود. واحدهای تولیدی همواره به دنبال حفظ سهم بازار و کسب سود بیشتر هستند و تلاش می‌کنند با استفاده از بازخورد پیشنهادات قبلی، پیشنهادات جدید را به گونه‌ی تهیه کنند که سود خود را بیشینه کنند. در ضمن نهاد قانونگذار و بهره‌بردار نیز با وضع قوانین و اجراء صحیح مکانیزم‌های حراج و تسویه، درصدد حفظ ثبات و منافع خرد و کلان تمام ذینفعان بازار هستند. معمولاً تمام اهداف در یک راستا نمی‌باشند و با توجه به هزینه زیاد سرمایه‌گذاری، تصمیم نادرست یا وضع قانون نادرست ممکن است سبب خسارت‌های فراوان شود. بنابراین، مدل‌سازی و شبیه‌سازی واکنش‌های بازار در تصمیم‌گیری‌ها از اهمیت بالایی برخوردار است و با توجه به پیچیدگی‌های بازار، کاربرد روش‌های محاسبات هوشمند برای تحلیل و پیش‌بینی اطلاعات بازار جایگاه ویژه و مهمی به خود اختصاص داده است (اصغری اسکویی 1394).

مدل عامل‌محور با ارائه یک مدل محاسباتی و در نظر گرفتن رفتار عوامل و نحوه تعامل آنها امکان بازنمایی و شبیه‌سازی عملکرد بازار عمده‌فروشی برق را فراهم می‌کند. با طراحی و بکارگیری مدل‌های عامل‌محور، در واقع می‌توان تاثیر استراتژی‌های مختلف تصمیم‌سازی در واحدهای تولیدی اعم از برنامه‌تولید و سرمایه‌گذاری و همچنین تاثیر وضع قوانین جدید توسط نهادها را پیش از اجراء مورد بررسی و ارزیابی قرارداد. یکی از چالش‌های مهم واحدهای تولیدی، انتخاب استراتژی صحیح پیشنهاد قیمت یا برنامه‌تولید است. پیشنهاد قیمت یا برنامه‌تولید به صورت مستمر و روزانه ارائه می‌شود و واحدهای تولیدی لازم است با توجه به بازخورد بازار درباره ادامه استراتژی یا تغییر آن تصمیم‌گیری کنند. این مقاله به بررسی کاربرد یادگیری تقویتی در ارائه پیشنهاد و مقایسه عملکرد آن با استراتژی‌های دیگر در میزان سود واحدهای تولیدی می‌پردازد.

یادگیری تقویتی<sup>1</sup> از شاخه‌های یادگیری ماشین و محاسبات هوشمند بوده و یک یادگیری مبتنی بر سعی و خطا است که از روانشناسی رفتارگرایی الهام می‌گیرد. یادگیری

کاربرد یادگیری تقویتی در یک مدل‌سازی... 5

تقویتی در اقتصاد و نظریه‌بازی‌ها بیشتر به بررسی تعادل‌های ایجاد شده تحت عقلانیت محدود عامل‌ها می‌پردازد. در یادگیری ماشین با توجه به این که بسیاری از الگوریتم‌های یادگیری تقویتی از تکنیک‌های برنامه‌نویسی پویا استفاده می‌کنند معمولاً مسئله تحت عنوان یک فرایند تصمیم‌گیری مارکف مدل می‌شود. این مقاله از یک مدل عامل‌محور شناخته شده و معتبر بازار برق ایران استفاده نموده و ضمن ارائه توصیف ریاضی از استراتژی تصمیم‌گیری مبتنی بر یادگیری تقویتی، تاثیر عملکرد آن را بر سود واحدهای تولیدی در طول زمان بررسی کرده و نتیجه را با دو استراتژی دیگر (استراتژی تصادفی و استراتژی حریصانه) مقایسه می‌کند. در این مطالعه هر بار اجراء مکانیزم حراج و تسویه یک بازی تکرار ایستا از نوع همکارانه کامل فرض شده است.

ساختار مقاله به شرح زیر است. بعد از مقدمه، در بخش دوم معرفی اجمالی از مدل‌سازی عامل‌محور از بازار برق ارائه می‌شود. بخش سوم پیشینه پژوهشی در حوزه کاربرد یادگیری تقویتی در مدل‌های عامل‌محور بازار برق ارائه می‌گردد. بخش چهارم یادگیری تقویتی برای مدل‌های تک‌عاملی، فرایند مارکفی و مدل‌های چندعاملی معرفی شده و سپس کارکرد این مدل در قالب تئوری بازی و به صورت یک بازی تکرار ایستا بیان می‌شود. در بخش پنجم، مدل بازار برق ایران و مفروضات مهم آن مختصراً شرح داده می‌شود. در بخش ششم نتایج پیاده‌سازی مدل و شبیه‌سازی استراتژی‌های مختلف بر عملکرد بازار ارائه شده و نتایج مورد بحث و بررسی قرار می‌گیرد. بخش هفتم، شامل جمع‌بندی و نتیجه‌گیری است.

## 2. مدل‌سازی عامل‌محور بازار برق

مطالعه در مورد نحوه تصمیم‌گیری و تحلیل عملکرد افراد و بنگاه‌ها در سامانه‌های ترکیبی و تعاملی، از قبیل یک جامعه یا یک بازار، همواره موضوع مورد توجه پژوهشگران علوم اجتماعی و حتی مهندسين بوده است. آنها با فرض عقلانی بودن عملکرد افراد و تبادل صحیح اطلاعات، تعامل اجتماعی را در حالت تعادل، بهینه تشخیص می‌دهند. مدل‌سازی

عامل محور (ABM) از این دیدگاه، عملکرد یک سامانه پیچیده اقتصادی را براساس رفتار و عملکرد تک تک عناصر تشکیل دهنده آن، توصیف و تحلیل می کند. این دیدگاه با توجه به محدودیت روش های تحلیل سنتی و پیشرفت در زمینه محاسبات، بیش از پیش مورد توجه قرار گرفته است. در مدل عامل محور، تصمیمات توسط یک ماهیت خودگردان و مستقل به نام عامل<sup>1</sup> گرفته می شود. این تصمیمات براساس مجموعه ادراکات عامل از محیط و دانش درونی عامل گرفته می شود. البته دانش و معرفت عامل در طول زمان با توجه به تجربه یا به صورت سعی و خطا تکامل یافته و ارتقاء می یابد. عامل به صورت عقلانی تصمیم می گیرد و قابلیت یادگیری، تطبیق و بازتولید دارد. تصمیم عقلانی، تصمیمی است که از تمامی اطلاعات در دسترس برای رسیدن به بیشینه سود استفاده می کند.

مدل عامل محور یک مدل محاسباتی است و از تعامل و تاثیر متقابل عامل ها بریکدیگر، عملکرد کل پدیده مورد نظر را شبیه سازی می نماید. تعامل و تاثیر متقابل عامل ها در سطح خرد، سبب بروز نتایج در سطح بازار می گردد. خصوصیت ویژه این نوع مدل، قابلیت تعامل رقابتی و مکرر عامل ها است. استفاده از مدل عامل محور مزایای متعدد دارد. از جمله اینکه امکان شبیه سازی و بازتولید عملکرد سامانه های پیچیده را به صورت محاسباتی فراهم می کند. خروجی یک سامانه ترکیبی، لزوما حاصل جمع مولفه های تشکیل دهنده نیست، بلکه برآیند تاثیر متقابل آنها بر یکدیگر است که مدل عامل محور امکان دستیابی به آن را میسر می سازد. مدل عامل محور امکان بازنمایی و توصیف ساده و طبیعی از یک سامانه ترکیبی و پیچیده را میسر می نماید. نهایتا اینکه، این نوع مدل سازی قابلیت انعطاف کافی برای نمایش سطوح مختلف از پیچیدگی اعم از تعداد و تنوع عامل ها، نحوه عملکرد عقلانی، عدم قطعیت، سیاست یادگیری و تکامل عامل ها را دارد.

---

1. Agent

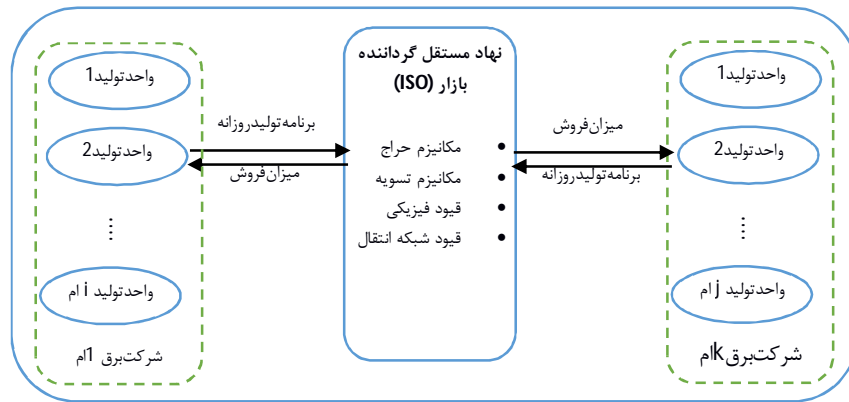
چنانچه اشاره شده مدل‌سازی عامل‌محور یک رویکرد طراحی پائین‌به‌بالا<sup>1</sup> است و در مسائل اقتصادی، که در آن‌ها نتایج بازار از برآیند تاثیر متقابل عوامل خرد شکل می‌یابد، کاربرد فراوان دارد. بازار عمده‌فروشی برق با مدل عامل‌محور به خوبی قابل بازنمایی است. بازیگران اصلی این بازار، واحدهای تولید برق هستند که به صورت مستقل و خودگردان تصمیم می‌گیرند به چه میزان و قیمت نیروی برق به بازار عرضه کنند. این واحدها با تبادل شفاف اطلاعات، امکان تعامل با دیگر بازیگران بازار را دارند و با توجه به معیار و قواعد درونی و اطلاعات بیرونی، استراتژی مناسب خود را با هدف بیشینه نمودن سود، انتخاب می‌نمایند. واحدهای تولیدی به لحاظ نوع فناوری و مشخصات فنی به چهار نوع اصلی نیروگاه‌های گازی، بخاری، ترکیبی و آبی‌هسته‌ای تقسیم می‌شوند. انواع عامل‌های دیگر نیز می‌تواند در این مدل وجود داشته باشد، که از جمله شرکت‌های انتقال و توزیع را می‌توان نام برد. در اغلب موارد برای ساده‌سازی مدل، نقش عامل‌های انتقال و توزیع به صورت قیود و قواعد محدودکننده مدل در نظر گرفته می‌شود. ضمناً نهاد بهره‌بردار (ISO) نیز وظیفه کنترل و مدیریت بازار را برعهده دارد. شکل (1) یک نمای از عوامل اصلی تشکیل‌دهنده مدل و جریان تبادل اطلاعات را نمایش می‌دهد. در مدل بازار برق، هر واحد تولیدی (عامل) با ارائه برنامه تولید<sup>2</sup> (SF) در مکانیزم حراج شرکت می‌کند. برنامه تولید شامل محدوده مشخص از میزان تولید (مگاوات) بازا قیمت‌های (ریال) مشخص است که برای یک دوره (24 ساعته) تهیه می‌شود. شکل (2) نمونه یک برنامه تولید را نشان می‌دهد. در این نمودار، قیمت پیشنهادی به ازاء میزان تولید به صورت پله‌ای در 10 سطح مشخص شده است. برنامه تولید پیشنهادی روزانه توسط هر عامل به نهاد بهره‌بردار ارائه می‌شود. بهره‌بردار با تجمیع برنامه تولید پیشنهادی واحدها و با توجه به پیش‌بینی تقاضای مصرف، در مکانیزم حراج قیمت بازار و فهرست واحدهای برنده و میزان خرید از هر یک را مشخص می‌کند. به عبارتی، با مکانیزم حراج، میزان فروش هر واحد که معیار سنجش موفقیت آن

---

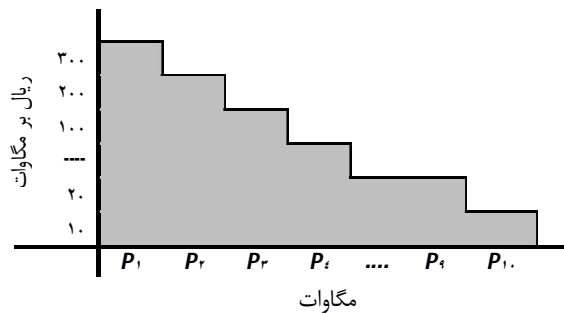
1. Bottom-Up

2. Supply Function (SF)

واحد یا پاداش<sup>1</sup> برنامه تولید (SF) است، مشخص می‌شود. معیار پاداش به عنوان بازخورد در تکرارهای بعدی جهت تصحیح یا حفظ استراتژی عامل مورد استفاده قرار می‌گیرد. بنابراین عامل ممکن است استراتژی تصمیم‌گیری خود را با دریافت بازخورد تغییر دهد.



شکل 1: نمائی از مدل عامل محور بازار عمده‌فروشی برق: واحدهای تولید برق به صورت عامل مستقل و در قالب گروه شرکت‌های برق می‌باشند و عامل‌های انتقال و توزیع به صورت قیود فیزیکی و فنی لحاظ شده است.



شکل 2: نمونه برنامه تولید (SF) پیشنهادی یک واحد به ازای یک ساعت مشخص



بهره‌بردار (ISO) وظیفه اجراء صحیح حداقل دو مکانیزم را دارد: حراج<sup>1</sup> و تسویه حساب<sup>2</sup>. از مکانیزم حراج برای تعیین قیمت بازار، لیست برندگان و میزان تولید هر عامل استفاده می‌شود. مکانیزم تسویه حساب برای تعیین میزان پرداخت به واحدها مورد استفاده قرار می‌گیرد. در حال حاضر از دو مکانیزم متفاوت حراج استفاده می‌شود: حداقل هزینه پیشنهادی<sup>3</sup> و حداقل هزینه پرداختی<sup>4</sup>. در اولی انتخاب تولیدکنندگان به گونه‌ی است که هزینه کل آنها حداقل گردد و در دومی به گونه‌ی انتخاب می‌شوند که هزینه‌های واقعی پرداخت مصرف‌کنندگان حداقل گردد. مکانیزم تسویه حساب نیز به دو روش قابل اجراء است. مکانیزم تسویه بر مبنای پیشنهاد<sup>5</sup> به این صورت است که به تولیدکنندگان منتخب همان مقداری که پیشنهاد داده‌اند، پرداخت می‌شود. در مکانیزم تسویه بر مبنای قیمت تسویه بازار<sup>6</sup> به کلیه تولیدکنندگان، قیمت یکسان که همان قیمت تسویه بازار است، پرداخت می‌شود.

چنانچه اشاره شد بکارگیری مدل عامل‌محور، فعالان عرصه بازار را قادر خواهد ساخت که فضای تصمیم‌گیری غیرمتمرکز و نامتجانس را شبیه‌سازی و واکاوی نمایند. عامل‌ها در بستر بازار و بر اساس اهداف، رویه تصمیم‌سازی، و الگوی رفتاری منحصر بفرد خود به کنش‌های متقابل خواهند پرداخت. عامل‌ها قادر خواهند بود که تجربیات گذشته خود را مبنای یادگیری ساخته، رفتار خود را با شرایط محیط تطبیق‌داده، و از فرصت‌های پدید آمده بهره‌جویند. از این‌رو مدل‌سازی عامل‌محور در زمره روش‌های تطبیقی و پیچیده قرار می‌گیرد. زنجیره این کنش و واکنش عامل‌ها، شبیه‌سازی بازار برق را پیش خواهد راند.

- 
1. Auction
  2. Settlement
  3. Offer Cost Minimization
  4. Payment Cost Minimization
  5. Pay-As-Bid (PAB)
  6. Market Clearance Price (MCP)

موفقیت و یا شکست عامل‌ها مبنای تطبیق راهبرد آن‌ها را شکل داده، و الگوی تکامل رفتار آن‌ها قابلیت تحلیل و بررسی خواهد یافت.

### 3. پیشینه پژوهش

راشدی و همکاران<sup>1</sup> (2016) یادگیری تقویتی را یک استراتژی پیشنهاد قیمت بهینه برای مدل عامل‌محور بازار عمده‌فروشی برق معرفی نمودند. آنها مدل چندعاملی با مکانیزم تسویه برمبنای قیمت تسویه بازار (MCP) و تعادل تابع عرضه<sup>2</sup> (SFE) را استفاده کرده و در قالب یک مسئله نظریه‌بازی با اطلاعات ناکامل و غیرهمکارانه، استراتژی بهینه را با کمک یادگیری تقویتی مشخص نمودند. در این مطالعه از الگوی فرایند تصمیم مارکوفی برای مدل تک‌عاملی استفاده شده و یک الگوی تصمیم برای مدل چندعاملی در بازی‌های تصادفی تکرارپذیر ارائه شده است. استراتژی پیشنهاد قیمت برای یک روز بعد بازار هر واحد تولیدکننده با کمک یادگیری تقویتی در مدل چندعاملی غیرهمکارانه و براساس استاندارد توزیع IEEE-30 پیاده سازی شده است. در این رویکرد یادگیری، برای استراتژی پیشنهاد قیمت، نیازی به تخمین حالت و فضای حالت نیست و هر عامل صرفاً برمبنای یادگیری از تجربه اختصاصی خود، پیشنهاد قیمت بعدی را ارائه می‌دهد.

راشدی و کبریائی<sup>3</sup> (2014) در یک مطالعه دیگر با استفاده از مفهوم تعادل نش، رفتار بازیگران بازار عمده‌فروشی برق را در دو حالت رقابتی و تعاملی مورد مطالعه قرار دادند. آنها برای این بازار، مدل بازی برمبنای تعادل تابع عرضه (SFE) را معرفی می‌نمایند. در این مدل تعدادی واحد تولیدکننده وجود دارد که هدف آنها پیشینه نمودن سود است. طبق تعریف تابع عرضه، سطح مطلوب تولید هر واحد را بازار قیمت‌های متفاوت مشخص

---

1. Rashedi, *et al.* (2016)

2. Supply function Equilibrium (SFE)

3. Rashedi and Kebriaei (2014)

می‌کند. در بازی رقابتی (ناهمکارانه)<sup>1</sup>، هر واحد صرفاً هدف بیشینه نمودن سود خود را دنبال می‌کند و نقطه تعادل شرایطی است که سود هر واحد در قبال پیشنهاد سایر واحدها، بیشینه باشد (تعادل نش). در بازی همکارانه، با توجه به این واقعیت که افزایش سود یک واحد سبب کاهش سود واحد دیگر می‌شود، عامل‌ها به دنبال سود بهینه همه‌جانبه<sup>2</sup> هستند و برای نیل به آن برخی قواعد منصفانه باید لحاظ شود. قوانین به نوعی است که به جای حداکثر نمودن سود هر عامل، یک تابع هدف تعریف می‌شود که سود مجموعه را نیز تامین نماید. بنابراین، تابع هدف از دو مولفه تشکیل شده است و علاوه بر سود هر عامل، میزان هماهنگی با سود مجموعه را نیز در بردارد. مولفه دوم ممکن است به صورت افزایشده یا کاهشده تاثیر گذارد. این مدل در دو حالت رقابتی و همکارانه با دو فرض تقاضای مصرف ثابت و تقاضای مصرف متناسب با قیمت، تحلیل شده و نتایج شبیه‌سازی آن مورد بررسی قرار می‌گیرد. این تحقیق به صورت تحلیلی اثبات می‌کند تعادل نش برای بازی رقابتی وجود دارد. در بازی همکارانه، در حالت تقاضای مصرف ثابت، بسته و محدب بودن تابع هدف شرط لازم و کافی برای وجود نقطه تعادل است. اما در حالت تقاضای مصرف متناسب با قیمت برق این شروط برای وجود نقطه تعادل کافی نمی‌باشند. رحیمیان و رجبی‌مشهدی<sup>3</sup> (2010) یک رویکرد تطبیقی یادگیری تقویتی برای بازار عمده‌فروشی برق که به صورت عامل‌محور مدل شده‌است، ارائه می‌دهند. در این تحقیق تلفیق دو رویکرد اکتشاف و بهره‌برداری<sup>4</sup> در یادگیری تقویتی بر اساس خصوصیات لحظه‌ای بازار برق به صورت تطبیقی تنظیم می‌شود. برای این منظور از یک سامانه استنتاج فازی برای تصمیم‌گیری و تنظیم پارامتر یادگیری بر اساس توان تولید هر عامل استفاده می‌شود. مزایای این روش، تطبیق پارامترهای یادگیری تقویتی بر اساس شرایط بازار و تولید هر

---

1. Non-cooperative

2. Pareto solution

3. Rahimiyan, *et al.* (2010)

4. Exploration and exploitation

واحد، کاهش حجم محاسبات مربوط به یادگیری، تعیین میزان ریسک هر استراتژی برای پیشنهاد قیمت و تدوین قوانین محاوره‌ی فازی برای انتخاب استراتژی است. بخشی از این تحقیق در رساله رحیمیان (1390) با عنوان رویکرد یادگیری تقویتی با خصوصیت فازی به عنوان یک استراتژی پیشنهاد قیمت در بازار عمده فروشی برق معرفی شده است. این رویکرد با روش قیمت گذاری متداول عامل محور مقایسه شده است.

ناظمی و همکاران (1390) میزان رقابت در بازار عمده‌فروشی برق براساس شاخص‌های ساختاری و نیز پتانسیل بروز رفتار غیررقابتی را مورد مطالعه قرار داده است. در این تحقیق مدلی از بازار برق ارائه شده و ضمن تعیین شاخص‌ها، ضریب همبستگی آنها مربوط به سال 1388 محاسبه شده است. سپس نتایج شبیه سازی مدل و تفاضل تعادل بهینه و رفتار راهبری محاسبه شده است. مدل این پژوهش علاوه بر قیود اقتصادی، قیود فنی را نیز شامل می‌باشد. موید کاظمی و شیخ‌الاسلامی (1393) با ارائه یک مدل عامل‌محور بازار سرمایه‌گذاری صنعت برق، استراتژی مبنی بر یادگیری تقویتی و تاثیر آن بر دینامیک سرمایه‌گذاری را بررسی و مطالعه نموده است. این تحقیق اثر راهبردهای تشویقی برای تامین هزینه‌های تولید و هزینه ثابت را بررسی می‌کند.

کراوس<sup>1</sup> و همکاران (2006) یک مدل عامل‌محور از بازار عمده‌فروشی برق را در نظر گرفته و ضمن استفاده از استراتژی یادگیری تقویتی شرایط رسیدن به تعادل نش را در صورت وجود، مورد بررسی قرار داده است. این تحقیق با استفاده از نتایج شبیه سازی مدل بازار نشان می‌دهد که در صورت وجود یک نقطه تعادل نش، عملکرد بازار بر اساس استراتژی یادگیری تقویتی به سمت آن همگرا است و در صورت وجود بیش از یک نقطه تعادل، بین آنها نوسان می‌کند. این تحقیق اثبات می‌کند که این همگرایی تحت پارامترهای مختلف بازار استوار می‌باشد. نتیجه مهم اینکه در صورت اثبات وجود نقاط تعادل مختلف، عملکرد سیستم نوسانی خواهد بود. مشیری و همکاران (1397) نیز از یک مدل عامل‌محور

---

1. Krause, et al. (2006)

کاربرد یادگیری تقویتی در یک مدل‌سازی... 13

برای شبیه‌سازی بازار برق ایران و ارزیابی سناریوی تغییر مکانیزم تامین سوخت نیروگاه‌ها و رقابتی کردن آن استفاده کردند و نشان دادند که تامین سوخت رقابتی در افزایش کارایی نیروگاه‌ها تاثیر مثبتی خواهد داشت.

#### 4. یادگیری تقویتی

یادگیری تقویتی، یکی از شاخه‌های یادگیری ماشین است که از روانشناسی رفتارگرایی الهام گرفته و بر رفتارهایی تمرکز دارد که عامل باید برای بیشینه کردن پاداش انجام دهد. طبق تعریف، رفتار عبارت از نگاهی است که از مجموعه ادراکات عامل به تصمیم‌سازی برای انجام اقدام مناسب منجر می‌شود و رفتار عقلانی، رفتاری است که با استفاده از تمامی اطلاعات در دسترس به دنبال بیشینه شدن کارآمدی می‌گردد. یادگیری تقویتی یک روش یادگیری برای انتخاب رفتار مناسب براساس پاداش و تنبیه است بدون اینکه لازم باشد نحوه انجام عمل را برای عامل مشخص نمائیم (حاج‌رسولیه‌ها 1393).

در یادگیری تقویتی، نوع اقدام عامل از قبل مشخص نمی‌شود، بلکه عامل با جستجوی مبتنی بر سعی و خطا رفتاری را یاد می‌گیرد که بیشترین پاداش را به دست آورده و سود کوتاه‌مدت فدای سود بلندمدت شود. در استراتژی جستجو برای رسیدن به پاداش بیشتر همواره دو رویکرد اصلی وجود دارد: رویکرد بهره‌مندانه<sup>1</sup> (حریصانه) و رویکرد اکتشافی<sup>2</sup> (تصادفی). چالش اصلی ایجاد یک تعادل با ترکیب دو رویکرد فوق است و لذا باید بین کاوش موارد جدید و استفاده از دانش پیشین تناسب ایجاد نمود. لذا در فضای جستجو با توزیع تصادفی یکنواخت باید چندین تکرار انجام داد تا امید ریاضی پاداش بلندمدت را حداکثر نمود. یادگیری تقویتی، در زمینه‌های گوناگونی کاربرد دارد. از جمله می‌توان از نظریه بازی‌ها، نظریه کنترل، تحقیق در عملیات، نظریه اطلاعات، سامانه چندعاملی، هوش ازدحامی، آمار، الگوریتم ژنتیک، بهینه‌سازی بر مبنای شبیه‌سازی نام برد. یادگیری تقویتی

---

1. Exploiting (Greedy)

2. Exploring (Random)

در اقتصاد و نظریه بازی بیشتر به بررسی تعادل‌های ایجاد شده تحت عقلانیت محدود عامل‌ها می‌پردازد.

اصولا یادگیری به سه نوع تقسیم‌بندی می‌شود: یادگیری با الگو، یادگیری بدون الگو و یادگیری تقویتی. در یادگیری با الگو، داده‌ها با برچسب، به صورت زوج مشترک، در فرایند یادگیری استفاده می‌شود و معمولا هدف یافتن یک الگو ساخت‌یافته بین داده‌ها است. در یادگیری بدون الگو، داده‌ها بدون برچسب و معمولا براساس شباهت بایکدیگر تحلیل می‌شوند. در یادگیری تقویتی، داده‌ها براساس تابع هزینه (پاداش یا تنبیه) تحلیل می‌شوند. یادگیری تقویتی با یادگیری با الگو دو تفاوت عمده دارد، نخست اینکه در آن زوج مشترک ورودی و خروجی در کار نیست و رفتارهای ناکارآمد نیز از بیرون اصلاح نمی‌شوند و دیگر آنکه تمرکز زیادی روی کارآمدی تعاملی وجود دارد که نیازمند رسیدن به یک تعادل بین اکتشاف‌های جدید و بهره‌برداری از دانش اندوخته شده دارد. (حاج‌رسولیه‌ها 1393)

از آنجا که یادگیری تقویتی براساس تعامل با محیط و عامل‌های دیگر شکل می‌گیرد، در ادامه به معرفی یادگیری تقویتی در دو گروه مسائل تک عاملی و چندعاملی اشاره می‌شود.

#### 1-4. یادگیری در مدل تک عاملی

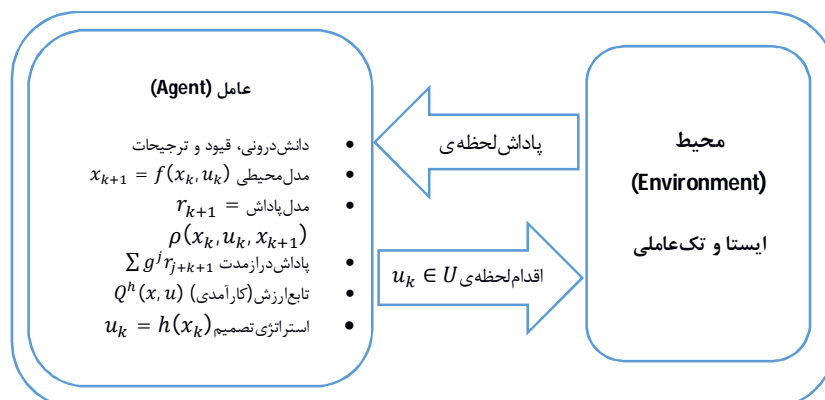
یادگیری در مدل تک‌عاملی، یک فرایند تصمیم‌مارکفی در نظر گرفته می‌شود. طبق تعریف در (باسینیو و همکاران<sup>1</sup> 2008)، فرایند مارکفی غیراحتمالی، شامل یک چهارتایی  $\{X, U, f, \rho\}$  است.  $X$  مجموعه فضای حالت و  $x_k \in X$  بردار حالت در لحظه  $k$  است.  $U$  مجموعه متناهی از اقدامات تعریف شده برای عامل است و اقدام عامل در لحظه  $k$  به صورت  $u_k \in U$  نمایش داده می‌شود. تابع  $f$  تابع تبدیلی است که تغییر حالت در اثر اقدام

---

1. Busoniu, et al. (2008)

کاربرد یادگیری تقویتی در یک مدل‌سازی... 15

عامل در لحظه  $k$  را به صورت  $x_{k+1} = f(x_k, u_k)$  توصیف می‌کند. به دنبال هر اقدام، عامل پاداش اسکالر  $r_k \in R$  را دریافت می‌کند که طبق تابع پاداش  $r_{k+1} = \rho(x_k, u_k, x_{k+1})$  محاسبه می‌شود. این پاداش صرفاً مربوط به حاصل اقدام در لحظه  $k$  است که محیط را از حالت  $k$  به حالت  $k+1$  منتقل می‌کند و ربطی به پاداش عملکرد درازمدت عامل ندارد. رفتار عامل براساس استراتژی<sup>1</sup> انتخاب اقدام مناسب به صورت  $u_k = h(x_k)$  مشخص می‌گردد. این استراتژی چنانچه نسبت به زمان ثابت باشد، ایستا نامیده می‌شود.



شکل 3: نمایی از تعامل مدل تک‌عاملی که به صورت فرایند مارکوفی غیراحتمالی توصیف می‌شود.

هدف یک عامل عقلانی، انتخاب اقدام مناسب در جهت بیشینه نمودن پاداش عملکرد در یک بازه درازمدت است. این معیار به صورت امید مجموع پاداش‌ها در بازه نامتناهی طبق رابطه (1) محاسبه می‌شود.

$$R_k = E \left\{ \sum_{j=0}^{\infty} \gamma^j r_{k+j+1} \right\} \quad (1)$$

در رابطه (1)  $r_k$  پاداش لحظه‌ای و  $R_k$  امید پاداش تجمعی است. ضریب  $\gamma$  عدد مثبت و کوچکتر از واحد است که ضریب تنزیل<sup>1</sup> نامیده می‌شود و میزان تاثیر پاداش لحظه‌ای در پاداش تجمعی را مشخص کرده و ضمناً سبب کراندار شدن آن می‌گردد. وظیفه عامل، انتخاب یک اقدام مناسب در لحظه  $k$  است به طوری که سبب بیشینه شدن پاداش تجمعی شود، لذا لازم است هر لحظه تخمینی از محاسبه پاداش تجمعی داشته باشد. برای این منظور تابع ارزش اقدام<sup>2</sup> که به اصطلاح تابع کیو<sup>3</sup> و به اختصار از این پس تابع ارزش نیز نامیده می‌شود، طبق رابطه (2) تعریف می‌گردد.

$$Q^h(x, u) = E \left\{ \sum_{j=0}^{\infty} \gamma^j r_{k+j+1} \mid x_k = x, u_k = u, h \right\} \quad (2)$$

تابع ارزش همواره وابسته به استراتژی تصمیم عامل ( $h$ ) است. شکل 3 نمائی از تعامل مدل تک‌عاملی و مولفه‌های اصلی را نشان می‌دهد. مقدار بهینه تابع ارزش، در بیشینه آن، به صورت  $Q^*(x, u) = \max_h Q^h(x, u)$  تعریف می‌شود. اثبات می‌شود که بیشینه تابع ارزش، طبق معادله بلمن<sup>4</sup> قابل محاسبه است (باسینیو و همکاران<sup>5</sup> 2008). براساس معادله بلمن، مقدار بهینه تابع، همواره معادل مجموع پاداش لحظه‌ای و مقدار بهینه تابع ارزش انتظاری در لحظه‌ی بعدی است. براین اساس، استراتژی بهینه، یک استراتژی غیراحتمالی و حریصانه<sup>6</sup> است که در هر لحظه، اقدامی را انتخاب می‌کند که بیشترین تابع ارزش را مطابق رابطه (3) محقق سازد.

$$h(x) = \operatorname{argmax}_h Q^h(x, u) \quad (3)$$

- 
1. Discount Factor
  2. Action-Value Function
  3. Q-function
  4. Bellman Optimality Equation
  5. Busoniu, et al. (2008)
  6. Greedy Policy



بنابراین یک عامل یادگیرنده، لازم است ابتدا تابع ارزش را هر لحظه تخمین زده و سپس براساس استراتژی حریمانه اقدامی را انتخاب کند که تابع ارزش را بیشینه نماید. لذا می‌توان گفت که یادگیری تقویتی، عبارت از تخمین تابع ارزش بهینه و یافتن استراتژی مناسب براساس آن است. تخمین تابع ارزش به دو صورت قابل انجام است: مبتنی بر مدل و بی‌نیاز از مدل. روش‌های برنامه‌ریزی پویا و اکتشافی (هیوریستیک)، یادگیری مبتنی بر مدل فضای حالت هستند. روش بی‌نیاز از مدل که به نام یادگیری کیو<sup>1</sup> نیز شناخته می‌شود، بسیار مفید و پرکاربرد است. یادگیری کیو یک روش محاسبات تکرارشونده است. در این روش، تخمین مقدار بهینه  $Q^*$  هر لحظه طبق رابطه (4) براساس پاداش لحظه‌ای بروزرسانی می‌شود.

$$Q_{k+1}(x_k, u_k) = Q_k(x_k, u_k) + \alpha_k [r_{k+1} + g \cdot \max_{u'} Q_k(x_{k+1}, u') - Q_k(x_k, u_k)] \quad (4)$$

$$Q_{k+1}(x_k, u_k) = (1 - \alpha_k) Q_k(x_k, u_k) + \alpha_k [r_{k+1} + g \cdot \max_{u'} Q_k(x_{k+1}, u')]$$

در رابطه (4) تخمین تابع ارزش، بدون نیاز به تابع تبدیل  $f$  و تابع پاداش  $\rho$  انجام می‌گیرد، لذا جزء روش‌های عاری از مدل شناخته می‌شود. ضریب یادگیری  $\alpha_k$  که یک عدد مثبت و کوچکتر از واحد است، نرخ بروزرسانی تابع ارزش بوده و سرعت همگرایی را تنظیم می‌کند. عبارت داخل کروشه، تفاضل تخمین  $Q_k(x_k, u_k)$  بازاء دو لحظه متوالی  $k$  و  $k+1$  را مشخص می‌کند که به اصطلاح نرخ تغییرات زمانی<sup>2</sup> نامیده می‌شود. اثبات می‌شود که مقدار  $Q$  در رابطه (4) تحت شرایطی به مقدار بهینه آن ( $Q^*$ ) همگرا می‌شود (باسینیو، 2008). یک شرط مهم همگرایی این است که تمام اقدام‌های ممکن برای حالت‌های محتمل (احتمال غیرصفر) لحاظ گردیده و جستجو برای یافتن مقدار بیشینه

---

1. Q-learning  
2. Temporal difference

تغییرات زمانی به صورت فراگیر باشد. بنابراین لازم است گاهی اوقات، اقدامی متفاوت از استراتژی حریمانه انتخاب شود که به آن استراتژی تصادفی یا اکتشافی گفته می‌شود. برای مثال، به ازاء هر اقدام، استراتژی حریمانه با احتمال  $\varepsilon$  و استراتژی تصادفی با احتمال  $1 - \varepsilon$  اختیار شود. رویکرد دیگر برای تامین شرط همگرایی، استراتژی اکتشافی بلتزن<sup>1</sup> است که طبق آن در حالت  $x$  اقدام  $u$  با احتمال رابطه (5) مشخص می‌شود.

$$h(x, u) = \frac{e^{Q(x,u)/\tau}}{\sum_{uw} e^{Q(x,w)/\tau}} \quad (5)$$

در رابطه (5) پارامتر  $\tau \in (0, \infty)$  که اصطلاحاً درجه حرارت نامیده می‌شود، میزان تصادفی یا حریمانه بودن استراتژی انتخاب را معین می‌کند. بازاء  $\tau \rightarrow 0$  استراتژی انتخاب کاملاً حریمانه و بازاء  $\tau \rightarrow \infty$  استراتژی انتخاب کاملاً تصادفی خواهد بود. استراتژی انتخاب، بازاء مابین این دو مقدار حدی، ترکیبی وابسته به پارامتر  $\tau$  خواهد بود. هرچه روند یادگیری پیشرفت کند، درجه حرارت بیشتر می‌گردد.

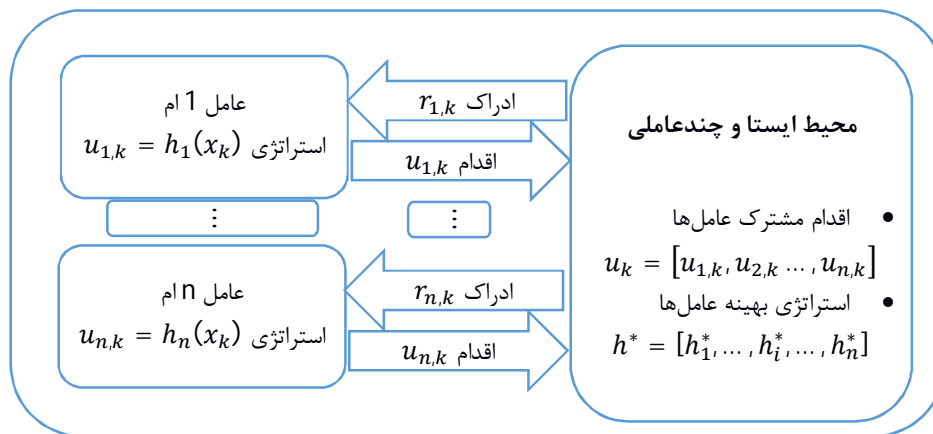
#### 2-4. یادگیری در مدل چند عاملی

حالت تعمیم یافته از فرایند تصمیم مارکوفی در مدل چندعاملی، تبدیل به بازی احتمالاتی<sup>2</sup> می‌گردد. در یک مدل مشتمل بر  $n$  عامل (شکل 4)، تغییر حالت از  $x_k$  به  $x_{k+1}$  به سبب اقدام مشترک<sup>3</sup> تمام عامل‌ها به صورت  $u_k = [u_{1,k}, u_{2,k}, \dots, u_{n,k}]$  تعریف می‌گردد. در این تعریف  $u_{i,k}$  اقدام عامل  $i$  ام در لحظه  $k$  ام است. اقدام مشترک و تغییر حالت، سبب پاداش  $r_{i,k+1} = \rho_i(x_k, u_{i,k})$  برای عامل  $i$  ام می‌شود. ضمناً هر عامل، استراتژی مخصوص به خود برای انتخاب اقدام دارد که به صورت  $u_{i,k} = h_i(x_k)$  نمایش داده می‌شود. در ادبیات نظریه بازی<sup>4</sup> چنانچه تابع پاداش تمام عامل‌ها یکسان باشد،

1. Boltzmann
2. Stochastic game
3. Joint action
4. Game theory

کاربرد یادگیری تقویتی در یک مدل‌سازی... 19

داشته باشد و  $\rho_1 = -\rho_2$  این بازی رقابتی<sup>2</sup> کامل یا مجموع صفر<sup>3</sup> نامیده می‌شود. بازی مخلوط<sup>4</sup> ترکیب‌های متنوعی است که نه همکارانه کامل و نه رقابتی کامل باشد. چنانچه ملاحظه می‌شود، فرایند تصمیم مارکفی و یادگیری تقویتی در مدل‌های چندعاملی وارد نظریه بازی‌ها شده و براساس ادبیات نظریه بازی قابل بیان و ترکیب می‌شود. از این پس برای اختصار و توصیف دقیق از اصطلاحات نظریه بازی استفاده خواهد شد (باسینیو و همکاران، 2008).



شکل 4: نمایی از مدل چندعاملی مشتمل بر n عامل.

### 3-4. بازی ایستا و تکراری

یک نوع خاص از بازی‌ها، مدلی است که مجموعه فضای حالت تهی است ( $X = \emptyset$ ) و پاداش صرفا وابسته به اقدام هر لحظه عامل دارد. به عبارتی پاداش صرفا تابع اقدام مشترک

1. Cooperative
2. Competitive
3. Zero-sum
4. Mixed games

عامل‌ها است  $r_{i,k+1} = \rho_i(u_k)$ . تهی بودن مجموعه فضای حالت به معنی این است که مدل، حافظه ندارد و تصمیم‌گیری بازاء هر اقدام (در محیط ایستا)، تکرار یک بازی مستقل است. البته با این تفاوت که در هر تکرار، از پاداش تکرارهای قبل در ساخت تصمیم بعدی می‌توان استفاده نمود. در این نوع بازی، از آنجا که فضای حالت وجود ندارد، استراتژی هر عامل تنها وابسته به مقدار ثابتی است که به بازی تخصیص داده شده است  $u_{i,k} = h_i(x = \text{constant})$ . عملکرد فرایند یادگیری تقویتی در این نوع بازی، معطوف به تعیین این مقدار ثابت در هر تکرار است، که مستقیماً در استراتژی و تبعاً در میزان پاداش لحظه‌ای و تجمعی تاثیر دارد.

در بازی ایستا یا تکراری، استراتژی تصادفی و چارچوب احتمالی، تقریباً اجتناب‌ناپذیر است. این موضوع به ویژه در توصیف ریاضی تعادل نش بسیار تعیین‌کننده است. فرض کنید در یک مدل چندعاملی، بهترین اقدام عامل  $i$  ام (که همزمان با اقدام سایر عامل‌ها است) و بیشترین پاداش انتظاری را به دنبال دارد با  $h_i^*$  مشخص شود، رابطه (6) برای هر تکرار بازی به صورت مستقل قابل کاربرد است (لذا اندیس  $k$  دیگر مورد استفاده قرار نمی‌گیرد).

$$E\{r_i|h_1, \dots, h_i, \dots, h_n\} \leq E\{r_i|h_1, \dots, h_i^*, \dots, h_n\} \forall h_i \quad (6)$$

تعادل نش<sup>1</sup> یک مجموعه استراتژی مشترک  $h^* = [h_1^*, \dots, h_i^*, \dots, h_n^*]$  است که انتخاب هر عامل بهترین انتخاب نسبت به سایر عامل‌ها باشد. تعادل نش، شرایطی را توصیف می‌کند که تمام عامل‌ها پاداش بیشینه ممکن را کسب نموده و تغییر استراتژی به سود هیچکدام نیست. اثبات شده است که هر بازی ایستا حداقل یک و بعضاً بیش از یک نقطه تعادل دارد. تعادل نش در اغلب برنامه‌های یادگیری تقویتی به عنوان هدف فرایند یادگیری در نظر گرفته می‌شود. لذا در شبیه‌سازی مدل چندعاملی، که عامل‌ها شروع به

---

1. Nash Equilibrium

کاربرد یادگیری تقویتی در یک مدل‌سازی... 21

تعامل با یکدیگر نموده و با توجه به الگوریتم یادگیری از عملکرد قبلی، اقدامی را انتخاب می‌کنند که پاداش بیشتری حاصل کنند، پس از گذشت زمان مشخصی به شرایطی می‌رسند که مجموع پاداش عامل‌ها ثابت شده و تغییر نمی‌کند، این حالت تعادل نش است (باسینیو، 2008).

#### 4-4. یادگیری مدل چندعاملی در حالت همکارانه کامل

در بازی همکارانه کامل، عامل‌ها تابع پاداش یکسان دارند ( $\rho_1 = \rho_2 = \dots = \rho_n$ ) و هدف یادگیری، بیشینه‌سازی پاداش تجمعی مشترک است ( $R = \sum R_{i,k}$ ). در این نوع بازی، عامل‌ها به دنبال سود بهینه همه‌جانبه هستند و شرایط به نوعی است که به جای صرفاً بیشینه نمودن سود هر عامل، یک تابع هدف کلی تعریف می‌شود که سود مجموعه را تامین نماید. با فرض وجود یک کنترل‌کننده مرکزی، این مسئله نیز به یک مسئله فرایند تصمیم مارکوفی در فضای مشترک استراتژی تبدیل می‌شود. در این حالت، یادگیری براساس تخمین تابع ارزش مشترک رابطه (7) و با استراتژی حریصانه، می‌تواند به پاسخ بهینه همگرا باشد. شروط همگرایی قبلاً اشاره شده است.

$$Q_{k+1}(x_k, u_k) = Q_k(x_k, u_k) + \alpha[r_{k+1} + g \cdot \max_{u'} Q_k(x_{k+1}, u') - Q_k(x_k, u_k)] \quad (7)$$

باتوجه به اینکه عامل‌ها مستقل هستند، هماهنگی آنها علی‌رغم یادگیری موازی و همزمان براساس یک تابع ارزش مشترک، هنوز یک چالش می‌تواند باشد. در استراتژی حریصانه، عامل باید به گونه‌ای اقدام کند که تابع ارزش طبق رابطه (8) حداکثر شود.

$$h_i^*(x) = \arg \max_{u_i} \max_{u_1 \dots u_{i-1} u_{i+1} \dots u_n} Q^*(x, u) \quad (8)$$

چنانچه اشاره شد، همگرایی وابسته به نوع استراتژی انتخاب عامل است. درعمل استراتژی حریصانه با مشکل مواجه شده و اغلب پاسخ نهایی به صورت بهینه محلی و نه ضرورتاً بهینه سراسری حاصل می‌شود. لذا در اینجا حل مسئله در سه حالت بدون

هماهنگی، با هماهنگی و هماهنگی غیرمستقیم، طرح می‌گردد. در رویکرد یادگیری بدون هماهنگی، ایده اصلی این است که لزوماً تابع ارزش بهینه، یکتا نیست. لذا در یادگیری تیمی، هر عامل براساس رابطه (7) تخمین جداگانه اختصاصی را داشته و براساس رابطه (8) استراتژی جداگانه را با هدف بیشینه نمودن پاداش، تعیین می‌نماید. الگوریتم یادگیری توزیعی<sup>1</sup> (لور و مارتین 2000) در حل مسائل همکارانه کامل در حالت غیراحتمالی به کار برده می‌شود. در این الگوریتم، هر عامل استراتژی اختصاصی  $h_i$  و تابع ارزش اختصاصی  $Q_i$  را دارد که صرفاً وابسته به اقدام همان عامل  $u_i$  است. ضمناً تابع ارزش اختصاصی، تنها زمانی بروزرسانی می‌شود که روند آن مطابق رابطه (9) افزایشی باشد.

$$Q_{i,k+1}(x_k, u_{i,k}) = \max\{Q_{i,k}(x_k, u_{i,k}), r_{k+1} + \max_{u_i} Q_{i,k}(x_{k+1}, u_i)\} \quad (9)$$

این شرط تضمین‌کننده آن است که تابع ارزش اختصاصی براساس رابطه (10) همواره حداکثر تابع ارزش مشترک را تقریب می‌کند.

$$Q_{i,k}(x, u_i) = \max_{u_1 \dots u_{i-1} u_{i+1} \dots u_n} Q_k(x, u) \quad \forall k, u = [u_1, \dots, u_n]^T \quad (10)$$

بنابراین مطابق رابطه (11) در یادگیری توزیع شده، استراتژی اختصاصی هر عامل فقط و فقط زمانی تغییر می‌کند که مستقیماً سبب ارتقاء تابع ارزش گردد.

$$h_{i,k+1}(x_k) = \begin{cases} u_{i,k} & \max_{u_i} Q_{i,k+1}(x_k, u_i) > \max_{u_i} Q_{i,k}(x_k, u_i) \\ h_{i,k}(x_k) & \text{otherwise} \end{cases} \quad (11)$$

به عبارتی، رویه رابطه (11) تضمین می‌کند که استراتژی اشتراکی  $h^* = [h_1^*, \dots, h_i^*, \dots, h_n^*]$  همواره براساس تابع ارزش سراسری  $Q^*$  بهینه خواهد بود. به

1. Distributed Q-Learning

2. Lauer and Martin (2000)

سادگی قابل اثبات است که به شرط مثبت بودن مقدار پاداش و تخصیص مقدار اولیه صفر به تابع ارزش اختصاصی هرعامل  $(\forall i, k Q_{i,0} = 0, r_{i,k} > 0)$  استراتژی اختصاصی عامل‌ها به سمت تابع ارزش مشترک بهینه همگرا خواهد شد (باسینیو و همکاران، 2008). در مدل بازار برق ایران که در قسمت بعد شرح داده شده است، شرایط فوق وجود دارد و استراتژی یادگیری توزیع شده به کار برده می‌شود.

## 5. مدل بازار برق ایران

هدف این مقاله بررسی تاثیر یادگیری تقویتی توزیع شده در عملکرد عامل‌های بازار برق ایران است. مدل نرم‌افزاری بازار برق ایران (فلاحی، 1391) که یک مدل چندعاملی است در پژوهشگاه نیرو پیاده‌سازی شده و سال‌هاست به صورت نرم‌افزار اصلی شبیه‌سازی و هدایت بازار مورد استفاده قرار گرفته است. مشخصه‌های این نرم‌افزار براساس داده‌های واقعی تنظیم شده و نتایج شبیه‌سازی بسیار نزدیک به عملکرد واقعی بازار برق ایران است (فلاحی، 1392). چنانچه در بخش‌های قبل توضیح داده شد، قیمت‌گذاری در بازار عمده‌فروشی برق، بازاری هرساعت در شبانه‌روز، توسط نهاد بهره‌بردار (ISO) بازار انجام می‌شود. با توجه به مفروضات این مدل، فرایند تعیین قیمت را می‌توان یک بازی ایستا فرض نمود که هرساعت تکرار می‌شود. در این بازی هرعامل قیمت پیشنهادی خود را به طور مستقل اعلام نموده و بهره‌بردار با توجه به تقاضای بار مصرفی، بهترین پیشنهادها را انتخاب می‌نماید. لازم به توضیح است که فرایند یادگیری سبب تغییر در استراتژی تصمیم‌سازی عامل می‌گردد، لکن روی مکانیزم بازار (محیط) که مستقل از عامل است، تاثیری از لحاظ پویایی ندارد. لذا ایستا و تکراری بودن صرفاً مربوط به محیط و مکانیزم حراج و تسویه است. یادگیری سبب پویایی رفتار عامل می‌شود، اما مکانیزم‌های اجرائی بازار را تغییر نمی‌دهد.

دراکثر سیستم‌های بهم پیوسته قدرت، گرچه بیشتر توان مورد نیاز توسط واحدهای حرارتی تامین می‌شود، طیف متنوعی از واحدها با فناوری و ظرفیت‌های مختلف در شبکه

موجود است. در مدل مورد مطالعه ما واحدهای از نوع بخاری، گازی، سیکل ترکیبی، آبی، هسته و دیزلی وجود دارد، که در این تحقیق تنها چهار نوع اصلی مورد بررسی قرار گرفته است. استراتژی‌های عملیاتی مختلفی برای انتخاب ترکیب بهینه از واحدها در جهت برآوردن تقاضا وجود دارد، که در ساعات مختلف روز نیز متغیر است. البته اغلب استراتژی مبتنی بر معیارهای اقتصادی مقدم بر معیارهای دیگر است. به بیان دیگر، معیار مهمی که در بهره‌برداری سیستم‌های قدرت وجود دارد، برآورده شدن تقاضای مصرف با حداقل هزینه و با بهره‌گیری از ترکیب بهینه نیروگاه‌های مختلف می‌باشد. به فرآیند تعیین زمان‌بندی و برنامه‌ریزی ورود و خروج واحدهای تولیدی اصطلاحاً مشارکت واحدها<sup>1</sup> (UC) اطلاق می‌شود. هدف اصلی UC حداقل کردن کل هزینه عملیاتی تولید است در حالی که همه محدودیت‌های متناظر با این مسئله نیز برآورده شوند. محدودیت‌ها به دو دسته تقسیم می‌شوند: (1) محدودیت‌های تولیدکنندگان، نظیر نرخ شیب و حداقل زمان روشن و خاموش بودن واحد و ... (2) محدودیت‌های سیستمی، نظیر میزان انرژی و ذخیره مورد نیاز سیستم، محدودیت شبکه انتقال و .... در تابع هدف نیز هزینه‌های مرتبط با تولید انرژی، شیب، هزینه‌های راه‌اندازی و خاموش کردن واحد به همراه تاثیر این تصمیمات روی درآمد یا هزینه‌های مشتریان نیز باید در نظر گرفته شوند. با توجه به تابع هدف و محدودیت‌های مطرح شده، UC یک مساله برنامه‌ریزی عدد صحیح مختلط غیرخطی با بزرگ مقیاس است (فلاحی، 1392).

نهاد بهره‌بردار (ISO) به طور کلی دو وظیفه مهم برای اجرای بازار بر عهده دارد: مکانیزم حراج و مکانیزم تسویه حساب. از مکانیزم حراج برای انتخاب پیشنهادات و همچنین تعیین میزان انرژی و خدمات جانبی استفاده می‌شود. سپس مکانیزم تسویه حساب برای تعیین میزان پرداخت به تولیدکنندگان منتخب مورد استفاده قرار می‌گیرد. در حال حاضر از دو مکانیزم حراج متداول در بازارهای عمده فروشی استفاده می‌شود: (1) مکانیزم

---

1. Unit Commitment (UC)



کمینه هزینه پیشنهادی<sup>1</sup> که در این روش انتخاب تولیدکنندگان به گونه‌ی است که هزینه کل آنها حداقل گردد. (2) مکانیزم کمینه هزینه پرداخت<sup>2</sup> که در این روش تولیدکنندگان به گونه‌ی انتخاب می‌شوند که هزینه‌های واقعی پرداخت مصرف‌کنندگان حداقل گردد. در مرحله دوم، بازارها با دو روش تسویه می‌شوند: (1) روش پرداخت بر مبنای پیشنهاد<sup>3</sup> (PAB) که در این روش به تولیدکنندگان منتخب همان مقداری که خودشان پیشنهاد داده‌اند پرداخت می‌شود. (2) روش پرداخت بر مبنای قیمت تسویه بازار<sup>4</sup> (MCP) که در این روش به تولیدکنندگان منتخب، قیمت یکسانی که همان قیمت تسویه بازار است پرداخت می‌شود. مسئله برنامه‌ریزی ورود و خروج واحدها (UC) و توزیع اقتصادی بار<sup>5</sup> (ED) جزء مسائل برنامه‌ریزی عدد صحیح هستند که هدفشان حداقل کردن هزینه تولید با توجه به تقاضای مصرف، ظرفیت واحد، احتیاجات ذخیره و سایر محدودیت‌های واحدهاست. این مسائل عموماً در رده مسائل پیچیده NP-hard قرار می‌گیرند. اما چون جز مسائل جدایی‌پذیر هستند به شکل کارآمدی با استفاده از روش‌های تجزیه و برنامه‌ریزی عددی قابل حل هستند و در حال حاضر بسته‌های نرم‌افزاری پیشرفته‌ای برای حل این دسته از مسائل در دسترس است (فلاحی، 1392). توضیحات بیشتر در مورد تعریف محدودیت‌ها و روش‌های عددی مربوط به حل مسئله، خارج از موضوع اصلی این مقاله است لذا در ادامه تنها به معرفی عامل‌ها و فرایند یادگیری جهت ارائه پیشنهاد قیمت آنها می‌پردازیم.

در مدل بازار برق ایران، مجموعاً 550 واحد تولیدکننده داریم که بازار حضور دارند. مشخصات فنی واحدهای تولیدی در پنج نوع فناوری شامل نیروگاه‌های بخاری، گازی، سیکل ترکیبی، آبی و هسته‌ای با جزئیاتی از قبیل حداقل ظرفیت ( $P_{i,min}$ ) و حداکثر ظرفیت ( $P_{i,max}$ ) تولید، متوسط هزینه تولید در ساعت ( $AVC_i$ )، برنامه تعمیرات سالانه و سایر

- 
1. Offer Cost Minimization
  2. Payment Cost Minimization
  3. Pay-As-Offer or Pay-As-Bid
  4. Pay-at-MCP
  5. Economic Despatching

موارد در مدل تعریف شده‌اند. هرعامل یا واحد تولید کننده، حداکثر 10 گزینه برای پیشنهاد قیمت دارد. براساس قوانین، این 10 گزینه براساس ضریبی از متوسط هزینه تولید هر واحد طبق رابطه (12) محاسبه می‌شود. این ضریب بزرگتر از یک و کوچکتر از دو می‌باشد.

$$h_i(j) = \left(1 + \frac{2 \times j}{10}\right) AVC_{i,j}, j = 0, 1 \dots 8 \quad (12)$$

هرعامل حداقل سه استراتژی برای قیمت گذاری دارد: تصادفی، حریصانه و یادگیرنده. در استراتژی تصادفی، حرکت اکتشافی است و قیمت پیشنهادی عامل  $i$  ام بازاء تکرار  $k+1$  ام بازی به صورت تصادفی (13) مشخص می‌شود.

$$u_{i,k+1} = h_i[\text{uniform.random}(0,8)] \quad (13)$$

در استراتژی غیرتصادفی، حرکت جهت‌دار است، لذا ابتداء لازم است تابع ارزش توزیعی، به صورت اختصاصی بازاء هرعامل، طبق رابطه (14) تخمین زده شود.

$$Q_{i,k+1}(j) = Q_{i,k}(j) + a[r_{i,k+1} - Q_{i,k}(j)] \quad \forall j \quad (14)$$

در این رابطه، تابع ارزش اختصاصی  $Q_{i,k}(j) = Q_{i,k}[u_{i,k} = h_i(j)]$  و پاداش اختصاصی  $r_{i,k+1} = \rho_i(u_k)$  برای عامل  $i$  ام حاصل از قیمت گذاری در لحظه  $k$  ام است. پارامتر آلفا، که عددی مثبت و کمتر از واحد است، ضریب یادگیری است. چنانچه قبلا اشاره شد، ضریب یادگیری نرخ همگرایی تخمین تابع ارزش به مقدار بهینه را مشخص می‌کند. استراتژی حریصانه، قیمتی را انتخاب می‌کند که بیشینه تابع ارزش تخمینی را برای تکرار بعدی سبب شود. لذا قیمت پیشنهادی عامل  $i$  ام در تکرار  $k+1$  ام از رابطه (15)

کاربرد یادگیری تقویتی در یک مدل‌سازی... 27

محاسبه می‌شود. یادآور می‌شود در این نوع بازی، فضای حالت تهی است و مسئله از نوع بازی تکرار است.

$$u_{i,k+1} = h_i \left[ \operatorname{argmax}_j Q_{i,k+1}(j) \right] \quad (15)$$

در استراتژی یادگیرنده، ترکیبی از استراتژی تصادفی و حریصانه مطابق رابطه (16) برای قیمت گذاری انجام می‌شود. این ترکیب براساس معادله بلتزن براساس معیار درجه حرارت  $T$  تنظیم می‌شود.

$$u_{i,k+1} = \begin{cases} h_i \left[ \operatorname{argmax}_j Q_{i,k+1}(j) \right] & \text{if } \text{uniform.random}(0,1) \leq \tau_k \\ h_i [\text{uniform.random}(0,8)] & \text{otherwise} \end{cases} \quad (16)$$

در این تحقیق معیار درجه حرارت به صورت تجربی براساس رابطه (17) تعیین می‌شود.

$$\tau_k = \min \left( \theta, \frac{12k + 500}{12k + 5200} \right) \quad (17)$$

در رابطه فوق  $k$  شمارنده تکرار بازی و پارامتر  $\theta = 0/93$  سقف احتمال استراتژی تصادفی را مشخص می‌کند (فلاحی، 1392). در پایان این بخش و قبل از بررسی نتایج شبیه‌سازی، مجدداً تاکید می‌شود که مدل مورد استفاده در این پژوهش برگرفته از مدل ارائه شده توسط پژوهشگاه نیرو است و این مقاله اثر استراتژی‌های مختلف تصمیم‌سازی را مطالعه می‌کند. در این مدل اثر منافع متقابل مصرف کننده و تولید کننده، در قالب برخی قیود<sup>1</sup> لحاظ شده است که از جمله رابطه (12) که سقف قیمت پیشنهادی عامل‌ها را مشخص می‌کند، می‌توان نام برد. ضمناً عدم امکان تبانی بین تولیدکنندگان و برونزا بودن

---

1. Constraints

مصرف مورد تاکید می‌باشد. با لحاظ این مفروضات و قیود، این مدل به شکل بازی همکارانه توصیف می‌شود، که عامل‌ها در تابع پاداش مشترک هستند، اما در زمینه سهم بازار رقابت دارند. در بازی همکارانه، عامل‌ها به دنبال سود بهینه همه‌جانبه هستند و برای نیل به آن برخی قواعد باید لحاظ شود. قوانین به نوعی است که به جای صرفاً بیشینه نمودن سود هر عامل، یک تابع هدف تعریف می‌شود که سود مجموعه را نیز تامین نماید.

## 6. نتایج شبیه‌سازی مدل

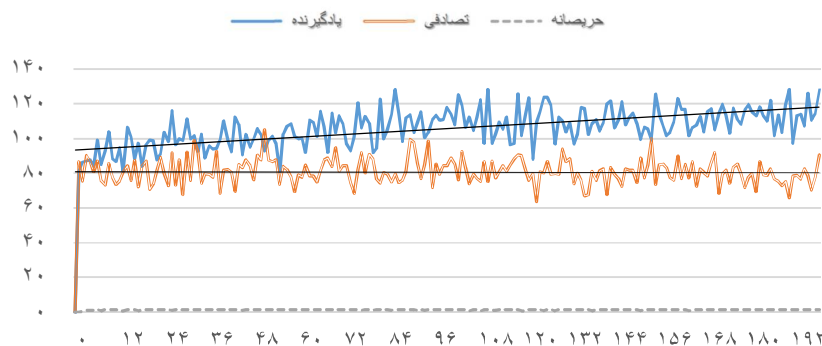
در این بخش نتایج شبیه‌سازی که برای مطالعه و بررسی تاثیر استراتژی یادگیرنده در عملکرد واحدها و سود عملیاتی آنها انجام شده است، ارائه می‌شود. شبیه‌سازی مدل بازار برق ایران، در حالت تسویه بر اساس قیمت پیشنهادی (PAB) و نحوه مشارکت واحدها (UC) در طول 200 تکرار با استراتژی‌های مختلفی انجام شده و نتایج ثبت شده است. با توجه به تفاوت ماهیت فنی و اقتصادی واحدها، ارزیابی‌ها به تفکیک فناوری نیروگاه‌ها انجام گرفته است. معیار ارزیابی عملکرد استراتژی‌ها، براساس پاداش یا سود فروش واحدها بازا هر مگاوات ثبت شده است.

### 6-1. تاثیر استراتژی یادگیری در سود واحدهای تولیدی

شکل‌های (5) تا (8) میانگین پاداش یا سود بازا هر مگاوات برای واحدهای تولیدی چهارفناوری اصلی، به ترتیب نیروگاه‌های بخاری، گازی، سیکل ترکیبی و آبی را نمایش می‌دهند. نتایج مربوط به 200 تکرار فرایند قیمت‌گذاری است که بر اساس سه استراتژی تصادفی، حریمانه و یادگیرنده عامل‌ها تصمیم‌گیری نموده و قیمت پیشنهادی را ارائه می‌دهند. از بررسی این چهار نمودار مشاهده می‌شود که در تمام واحدهای تولیدی، اعم از بخاری، گازی، ترکیبی و آبی، سود حاصل از استراتژی صرفاً حریمانه بسیار کمتر از دو استراتژی دیگر است. از لحاظ نظری اشاره شد که این روش در یک نقطه کمینه محلی متوقف شده و به پاسخ بهینه تقریب نمی‌کند. مقایسه سود واحدهای تولیدی نشان می‌دهد،

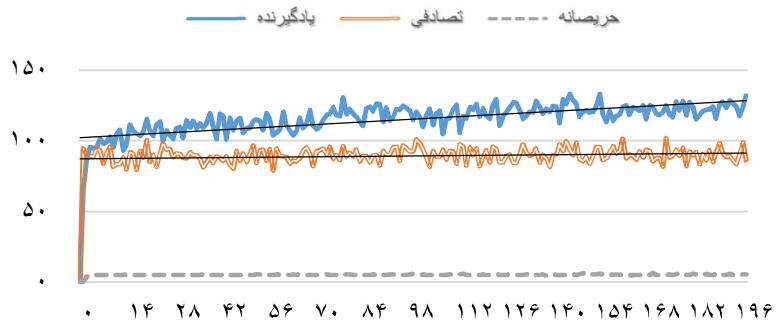
سود حاصل از استراتژی صرفاً تصادفی یک روند تقریباً ثابت دارد در حالیکه سود حاصل از استراتژی یادگیرنده یک روند افزایشی یکنواخت را دنبال می‌کند. شکل (9) نموداری از شیب تقریب خطی تغییرات سود را نمایش می‌دهد. چنانچه مشخص است، افزایش سود ناشی از استراتژی یادگیرنده در نیروگاه‌های حرارتی (بخاری، گازی و ترکیبی) بیش از ده درصد است که به صورت قابل توجهی بیشتر از استراتژی صرفاً تصادفی است. در نیروگاه‌های آبی افزایش سود با استراتژی یادگیرنده بیش از 2 درصد است که در مقایسه با استراتژی تصادفی بیشتر است.

اساساً سهم تولید نیروگاه‌های آبی در بازار ایران نسبت به نیروگاه‌های حرارتی کمتر است. شکل (10) سهم بازار فناوری‌های مختلف را نمایش می‌دهد. طبق نتایج شبیه‌سازی، سهم نیروگاه‌های بخاری، گازی و ترکیبی هر یک حدود 30 درصد است و مجموع سهم نیروگاه‌های آبی و هسته‌ای کمتر از 10 درصد است. شکل (11) نتایج تحلیل آماری مقایسه سود ناشی از استراتژی یادگیرنده با دو استراتژی دیگر را نشان می‌دهد. براساس این تحلیل توزیع پراکندگی پاداش واحدهای گازی بازا استراتژی یادگیرنده به صورت معناداری متمایز و بیشتر از استراتژی حریصانه، تصادفی است.

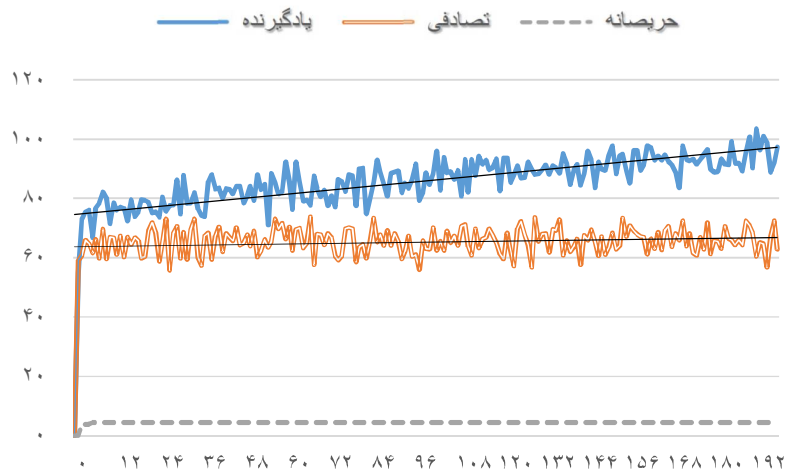


شکل 5: پاداش (سود) بازا هر مگاوات در واحدهای بخار براساس استراتژی حریصانه، تصادفی و یادگیرنده در

200 تکرار بازی

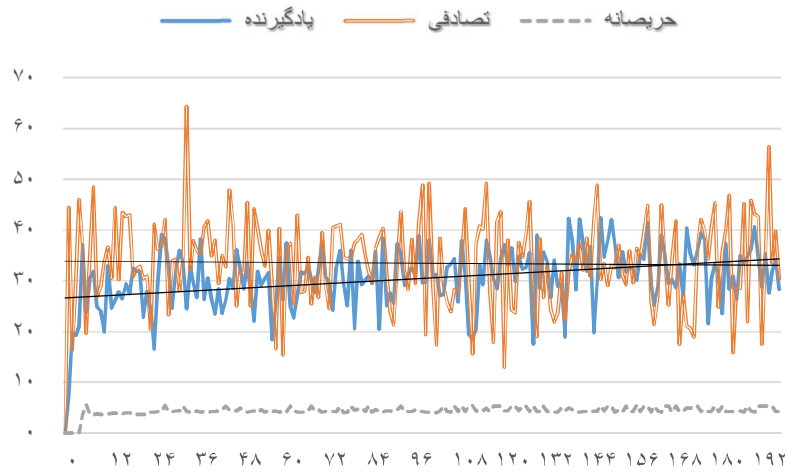


شکل 6: پاداش (سود) بازاء هرمگاوات در واحدهای گازی براساس استراتژی حریصانه، تصادفی و یادگیرنده در 200 تکرار بازی

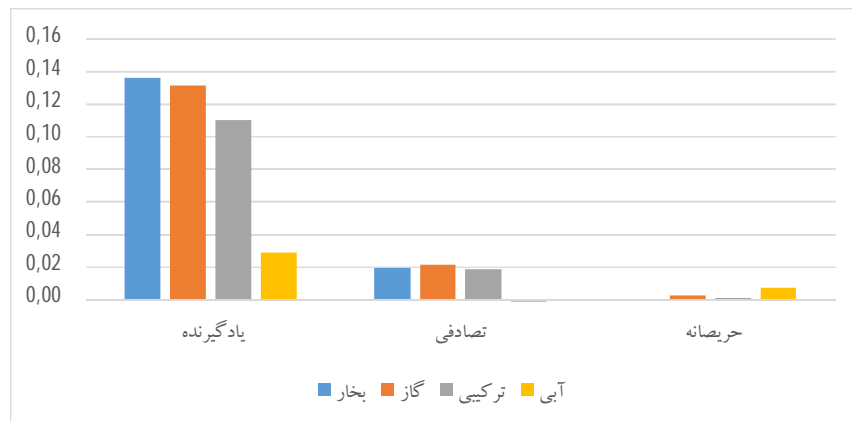


شکل 7: پاداش (سود) بازاء هرمگاوات و واحدهای سیکل ترکیبی براساس استراتژی حریصانه، تصادفی و یادگیرنده در 200 تکرار

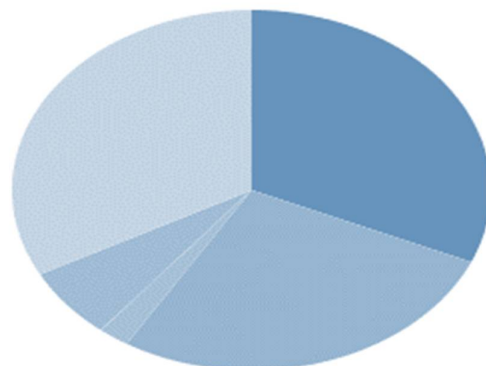
کاربرد یادگیری تقویتی در یک مدل‌سازی... 31



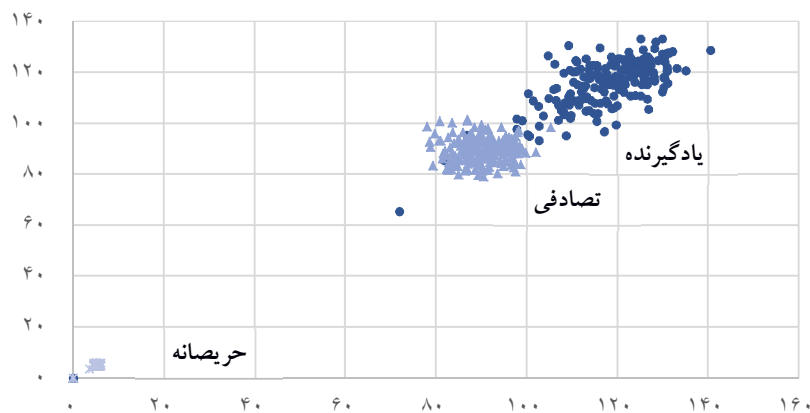
شکل 8: پاداش (سود) بازاء هر مگاوات واحدهای آبی براساس استراتژی حریمانه، تصادفی و یادگیرنده در 200 تکرار بازی



شکل 9: شیب تغییرات خطی پاداش (سود) واحدها براساس استراتژی حریمانه، تصادفی و یادگیرنده در طول 200 تکرار بازی



شکل 10: سهم بازار فناوری‌های مختلف در شبیه‌سازی بازار برق ایران

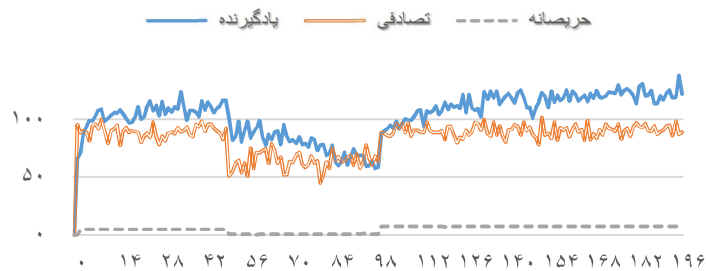


شکل 11: پراکندگی پاداش (سود) براساس استراتژی حریصانه، تصادفی و یادگیرنده در 200 بازی

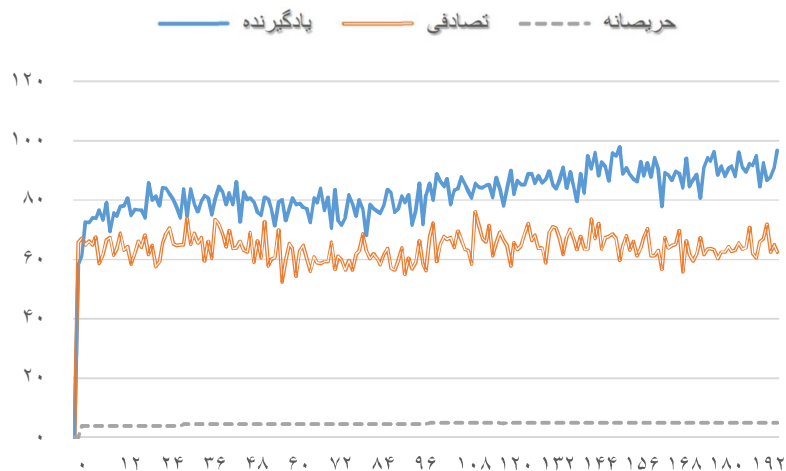


## 2-6. تاثیر شوک تغییر تقاضای بار مصرفی

آزمون دوم، با هدف بررسی عملکرد عامل‌ها در دوره گذرا و روند رسیدن به وضع پایدار طراحی گردید. در واقع عملکرد بازار تا رسیدن به نقطه تعادل در این آزمون مورد بررسی قرار می‌گیرد. این آزمون نیز شامل 200 بار تکرار بازی قیمت‌گذاری است، با این تفاوت که در تکرار پنجاهام تقاضای بار مصرفی ناگهان به سطح 80٪ کاهش داده می‌شود و مجدداً در تکرار صدم به میزان بار کامل بازگردانده می‌شود. این تغییرات بازار سه استراتژی مورد نظر، جداگانه مورد بررسی قرار گرفته و نتایج ثبت شده است. تبعاً تغییر بار، تاثیر مستقیم در پیشنهاد قیمت‌ها خواهد داشت. در این آزمون واکنش واحدها و سود آنها به تفکیک فناوری مورد بررسی قرار گرفته است. شکل (12) روند تغییر سود واحدهای گازی را بازار سه استراتژی حریصانه، تصادفی و یادگیرنده نشان می‌دهد. مشخص است استراتژی یادگیرنده چه در حالت بار کامل و چه در حالت کاهش بار ناگهانی، سود بیشتری برای واحدها تامین می‌نماید. البته واحدهای گازی به سبب ماهیت فناوری، اینرسی کمتری نسبت به تغییرات دارند. شکل (13) واکنش نیروگاه‌های سیکل ترکیبی را به این تغییرات نشان می‌دهد. مشخص است شوک بار، تاثیر محسوسی بر عملکرد و سود این نوع واحدها نداشته و البته طبق انتظار اولیه استراتژی یادگیرنده همواره قالب است. این موضوع در مورد نیروگاه‌های آبی نیز مشابه است.



شکل 12: پاداش (سود) واحدهای گازی با وجود شوک تغییر تقاضا به ازاء استراتژی حریصانه، تصادفی و یادگیرنده



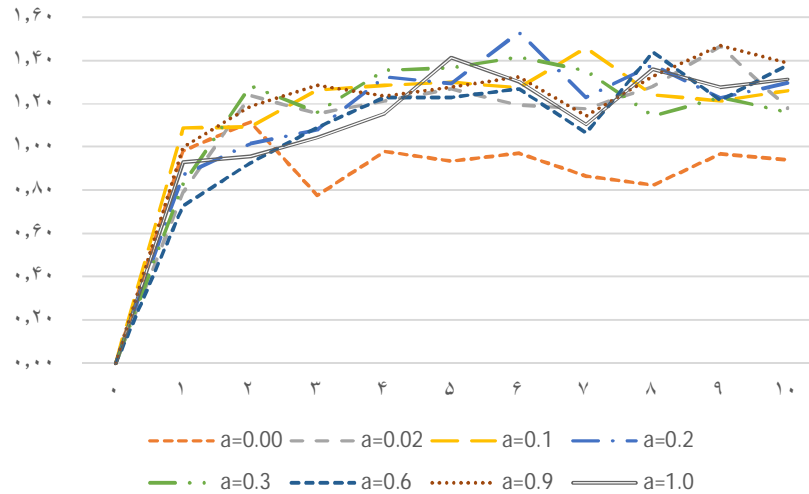
شکل 13: پاداش (سود) واحدهای سیکل ترکیبی با وجود شوک تغییر تقاضا به ازاء استراتژی حریصانه، تصادفی و پادگیرنده

### 3-6. تاثیر ضریب یادگیری

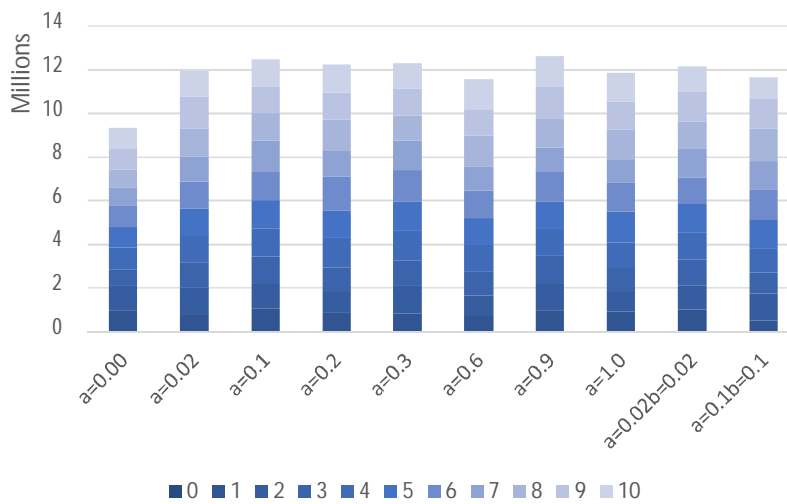
ضریب یادگیری، چنانچه شرح داده شد، در سرعت همگرایی و رسیدن به حالت تعادل در قبال تغییرات در سیستم، تاثیر دارد. این موضوع در شکل (14) مورد مطالعه قرار گرفته است. این شکل روند تغییر سود واحدهای بخار را در 10 تکرار اول به تفکیک مقادیر مختلف ضریب یادگیری نشان می دهد. چنانچه مشاهده می شود تفاوت معناداری بین نتایج با ضریب یادگیری صفر و غیرصفر دیده می شود، اما تفاوتی معنادار بین اثر مقادیر غیرصفر وجود ندارد. شکل (15) مجموع سود ثبت شده واحدهای بخار را در 10 تکرار اول به تفکیک ضرایب یادگیری نشان می دهد. در شکل اخیر شیب استراتژی تصادفی نیز به منظور مقایسه و درک تاثیر ضریب یادگیری، اضافه شده است. ضریب یادگیری در میزان سود یا پاداش اثری ندارد و تنها می تواند در سرعت همگرایی موثر باشد. جمع بندی این است که ضریب یادگیری غیرصفر که سبب تاثیرپذیری فرایند پیشنهاد قیمت از مقدار

کاربرد یادگیری تقویتی در یک مدل‌سازی... 35

پاداش در تکرارهای قبلی است، تاثیر به سزائی در افزایش سود واحدها دارد اما اندازه آن ضریب، تاثیر معناداری در میزان سود ندارد.



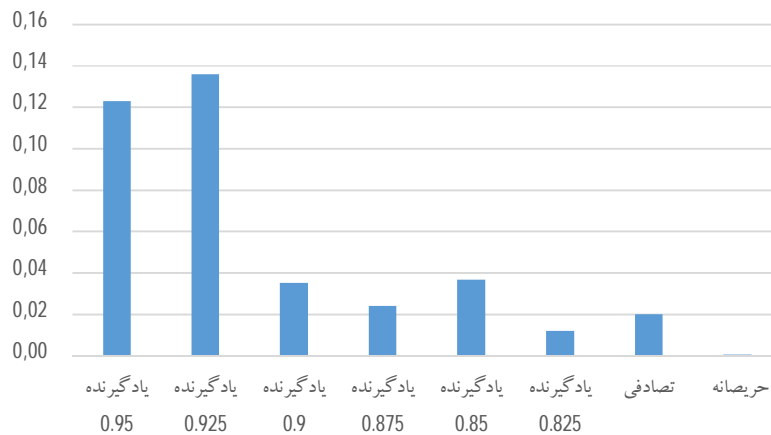
شکل 14: تغییرات پاداش در 10 تکرار اول بازی به ازاء ضرایب یادگیری متفاوت (واحدهای بخار)



شکل 15: مجموع پاداش واحدهای بخار در 10 تکرار اول بازماء مقادیر مختلف ضریب یادگیری

#### 4-6. تاثیر پارامتر درجه حرارت

چنانچه در بخش‌های نظری یادگیری شرح داده شد، معادله بهینه‌سازی بلتزن یک معیار نمادین به نام درجه حرارت معرفی می‌کند که میزان تلفیق دو استراتژی تصادفی و حریصانه را تعیین می‌کند. این معیار در مدل ما به صورت روابط (16) و (17) تعریف شده است. نام درجه حرارت از الگوریتم تبرید شبیه‌سازی<sup>1</sup> شده گرفته است که به تدریج که تعداد تکرار افزوده می‌شود از میزان استراتژی تصادفی کاسته شده و بر میزان استراتژی حریصانه افزوده می‌شود. شکل (16) اثر ناشی از ترکیب‌های مختلف از دو استراتژی تصادفی و حریصانه را در قالب شیب خطی سود عامل‌ها نمایش می‌دهد. این نمودار مشخص می‌کند که با دستکاری تغییرات درجه حرارت در طول مدت تکرار بازی، همگرایی تابع ارزش در استراتژی یادگیرنده مختل شده در نتیجه بر قیمت‌گذاری عامل‌ها و شیب افزایشی سود آنها به شدت تاثیر منفی می‌گذارد.



شکل 16: شیب تغییرات خطی پاداش واحدها براساس استراتژی یادگیرنده با درجه حرارت ( $\theta$ ) مختلف و مقایسه آن با استراتژی حریصانه، تصادفی

## 7. جمع‌بندی و نتیجه‌گیری

پیشرفت‌های فناوری دهه‌های اخیر، بازارهای برق را از حالت تمرکز به سوی رقابتی شدن هدایت کرده است. هرچند رقابتی شدن بازارها موجب کارآیی بیشتر و در نتیجه افزایش رفاه جامعه می‌شود، اما با توجه به ویژگی‌ها و پیچیدگی‌های خاص بازار برق، چالش‌های نوینی نیز به همراه دارد. بازارهای برق سنتی معمولاً توسط تولیدکنندگان انحصاری و تحت نظارت دولت تنظیم می‌شوند، اما میزان تولید و قیمت در بازارهای رقابتی از رقابت و تعامل تولیدکنندگان در یک بازار ایجاد شده توسط نهاد تنظیم بازار تعیین می‌شوند. با توجه به برقراری تعادل لحظه‌ای در بازار برق و عدم امکان ذخیره‌سازی برق در مقیاس وسیع، بازارهای برق به صورت یک روز یا یک ساعت قبل تنظیم می‌شوند. این بازار شامل تعداد زیادی تولیدکننده ناهمگن با فناوری‌ها و استراتژی‌های متفاوت و نهاد مستقل تنظیم بازار به عنوان حراج‌گر است که در چارچوب قواعد از پیش تعیین شده عمل می‌کنند. مدل‌سازی چنین بازاری با توجه به ناهمگونی عوامل مستقل تصمیم‌گیر و تعامل و یادگیری آنها در طول زمان با مدل‌سازی‌های کلاسیک که عوامل بازار را همگون فرض می‌کنند امکان‌پذیر نیست. مدل‌سازی عامل‌محور و شبیه‌سازی عملکرد بازار برق به عنوان یک رویکرد کارآمد برای مطالعه رفتار بازیگران بازار و تاثیر تعامل آنها بر یکدیگر و بر تعادل بازار در سال‌های اخیر مطرح شده است. این مدل‌ها توانایی تبیین رفتار عوامل ناهمگون و تعامل آنها را در یک محیط پویا که منجر به تعادل بازار می‌شود را دارند، اما به دلیل گستردگی و پیچیدگی آنها از شبیه‌سازی و تحلیل سناریو به جای حل معادلات سیستم استفاده می‌شود. یکی از چالش‌های مهم در بازارهای برق رقابتی، مطالعه رفتار عامل‌های تولیدی ناهمگن و استراتژی پیشنهاد قیمت آنها برای یک روز بعد است. هر عامل عقلانی با هدف حداکثر نمودن سود درازمدت، به دنبال انتخاب استراتژی مناسب برای پیشنهاد قیمت در مکانیزم حراج است. در این تحقیق سه استراتژی حریصانه، تصادفی و یادگیرنده مورد مطالعه و بررسی قرار گرفته است. مدل مورد استفاده یک مدل چندعاملی از بازار برق ایران است که مکانیزم حراج در آن یک بازی ایستا و تکراری در نظر گرفته شده است. فرایند

یادگیری تقویتی توزیع شده برای یادگیری عامل از نتیجه بازی های قبلی مورد استفاده قرار گرفته است. ارزیابی ها در حالت تعادل صورت گرفته و نتایج نشان می دهد سود حاصل از یادگیری تقویتی، بصورت معناداری بیش از استراتژی حریصانه و تصادفی است. در حالت شوک و تغییرات بار مصرفی، یادگیری تقویتی باعث همگرایی و تعادل شبکه می گردد. از محدودیت های این مطالعه می توان به عدم امکان تشکیل ائتلاف و تبانی در بازار و در نظرنگرفتن حالت بازی همکارانه هماهنگ شده در مدل و همچنین برونزا فرض کردن میزان مصرف روز بعد اشاره کرد که هر یک از آن ها می توانند موضوعات تحقیقات آتی قرار گیرند.

## 8. منابع

### الف) فارسی

اصغری اسکویی، محمدرضا (1394)، پیش بینی سری های زمانی مالی با کمک شبکه عصبی تاخیری، فصلنامه پژوهش های اقتصادی ایران، سال 15، شماره 57، صص 75 – 108.

رحیمیان، مرتضی (1390)، سنجش و تحلیل بازار در بازار برق به کمک اقتصاد محاسباتی عامل محور، رساله دکتری، دانشگاه فردوسی مشهد.

ناظمی، علی، خوش اخلاق، رحمان، عمادزاده، مصطفی، شریفی، علی مراد (1390)، برآورد قدرت بازار در بازار برق عمده فروشی ایران، تحقیقات مدل سازی اقتصادی، دوره 1، شماره 4، صص 31 – 55.

موید کاظمی، حمیدرضا و شیخ الاسلامی، محمد کاظم (1393)، ارزیابی راهبردهای توسعه تولید از دید رگولاتور با در نظر گرفتن دینامیک سرمایه گذاری بازیگران بازار، بیست و هشتمین کنفرانس بین المللی برق، تهران، ایران، صص 1-7.

کاربرد یادگیری تقویتی در یک مدل‌سازی... 39

فلاحی، فرهاد (1391)، «گزارش تحلیلی بازار برق ایران: کاربست شبیه‌سازی عامل محور در مطالعات بازار برق»، گزارش‌های پژوهشکده انرژی و محیط زیست پژوهشگاه نیرو، گروه پژوهشی اقتصاد و مدیریت برق.

فلاحی، فرهاد (1392)، «بررسی و تعیین مدل خرید انرژی الکتریکی در بازار روزانه از دید بهره‌بردار مستقل سیستم»، گزارش‌های پژوهشکده انرژی و محیط زیست پژوهشگاه نیرو، گروه پژوهشی اقتصاد و مدیریت برق.

مشیری، سعید، مروت، حبیب، ونصیری، عباس (1397)، بررسی تاثیر افزایش قیمت سوخت بر قیمت برق با استفاده از مدل‌سازی عامل بنیان بازار برق، فصلنامه مطالعات اقتصاد انرژی، سال 14، شماره 56، صص 12-22.

حاج‌رسولیه‌ها، حسین (1393). هوش مصنوعی پیشرفته، تالیف راسل و نورینگ، جلد دوم، چاپ سوم، تهران: نشر نیازدانش.

(ب) انگلیسی

Rashedi, N., Tajeddini, M. A., and Kebriae, H. (2016), Markov Game Approach for Multi-Agent Competitive Bidding Strategies in Electricity Market, *IET Generation, Transmission & Distribution*, Vol.10, No.15, pp. 3756-3763.

Rashedi, N., and Kebriaei, H. (2014), Cooperative and Non-Cooperative Nash Solution for Linear Supply Function Equilibrium Game. *Applied Mathematics and Computation*, Vol. 244, pp.794-808.

Rahimiyan, M., and Rajabi Mashhadi, H. (2010), An Adaptive Q-Learning Algorithm Developed for Agent-Based Computational Modeling of Electricity Market. *IEEE in Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, Vol.40, Issue.5, pp. 547-556.

Krause, Th., Vdovina Beck, E., Cherkaoui, R., Germond, A., Andersson, G. and Ernst, D. (2006), A Comparison of Nash Equilibria Analysis and Agent-Based Modeling for Power Markets. *International Journal of Electrical Power & Energy Systems*, Vol.28, Issue.9, pp.599-607.

Busoniu, L., Babuska, R. and De Schutter B.(2008), A Comprehensive Survey of Multiagent Reinforcement Learning. *IEEE in Transactions on Systems, Man, And Cybernetics-Part C: Applications and Reviews*, Vol.38, Issue.2, pp.1-18.

Lauer, M., and Riedmiller, M. (2000), an Algorithm for Distributed Reinforcement Learning in Cooperative Multi-Agent Systems. In

Proceedings of the Seventeenth International Conference on Machine Learning.

Toman, M. and Barbora, J. (2003), Energy and Economic Development: an Assessment of the State of Knowledge, Discussion Paper 03-13, Resources for the Future, Washington.