

به نام خدا



دانشگاه تهران
پردیس دانشکده‌های فنی
دانشکده مکانیک



کنترل آونگ معکوس دورانی به کمک الگوریتم‌های یادگیری تقویتی

پروژه تخصصی کارشناسی

رشته مهندسی مکانیک

نویسنده: آرش هاتفی

استاد راهنما: دکتر مسعود شریعت پناهی

بهمن ۱۳۹۹

تشکر و قدردانی

بدینوسیله، از آقایان جناب دکتر شریعت پناهی و جناب دکتر کاشانی برای حمایت‌های دلسوزانه‌شان در تمامی مراحل انجام این پروژه تشکر و سپاسگزاری می‌نمایم.

چکیده

در سال‌های اخیر، در پی پیشرفت‌های چشمگیر در هوش مصنوعی و یادگیری ماشینی، روش‌های هوشمند در حوزه‌ی کنترل محبوبیت فراوانی یافته‌اند. تاکنون عملکرد این الگوریتم‌ها همواره توسط بسیاری از مسائل مختلف به چالش کشیده شده و بسیاری از پیشرفت‌های صورت گرفته در روش‌های مذکور ضمن تلاش برای بهبود عملکردشان روی یک مسئله‌ی خاص اتفاق افتاده است. در این پروژه قصد داریم به بررسی عملکرد الگوریتم‌های یادگیری تقویتی، دسته‌ای از الگوریتم‌های کنترلی هوشمند، در کنترل آونگ معکوس دورانی (Rotary Inverted Pendulum) بپردازیم. دینامیک پیچیده و بسیار ناپایدار آونگ معکوس دورانی، همواره این سیستم را برای پژوهشگران حوزه‌ی کنترل جذاب نموده تا جایی که این مسئله را به‌عنوان یکی از مسائل معیار در حوزه‌ی مهندسی کنترل جهت مقایسه‌ی عملکرد روش‌های مختلف بدل نموده است. در گام نخست پروژه، یک آونگ معکوس دورانی جهت اجزای فرایند آموزش الگوریتم‌های هوشمند ساخته و آماده‌سازی شد. به موازات فرایند ساخت، به پیاده‌سازی یک مدل کامپیوتری دقیق از آونگ جهت تسریع فرایند آموزش در محیط شبیه‌سازی مجازی نیز پرداخته شد. در نهایت، سعی شد تا با بهره‌برداری مناسب از امکانات محاسباتی موجود، عملکرد دو الگوریتم یادگیری تقویتی از نوع گسسته (Deep Q-Learning) و پیوسته (Deep Deterministic Policy Gradient) بر روی مسئله‌ی کنترل آونگ معکوس دورانی سنجیده شود. هدف نهایی پروژه، دستیابی به درکی دقیق از نقاط ضعف و قوت روش‌های یادگیری تقویتی در مواجهه با مسائل پیچیده‌ی کنترل بی‌درنگ می‌باشد چراکه این کار گام اولیه و اساسی در راستای برطرف شدن نقاط ضعف و ایرادات وارد بر این روش‌ها است. عملکرد الگوریتم یادگیری تقویتی پیوسته، در محیط شبیه‌سازی از عملکرد الگوریتم گسسته پیشی گرفت. در نقطه‌ی مقابل، به علت کم‌تر بودن حجم محاسبات، الگوریتم گسسته توانست در عمل سازه‌ی واقعی را بهتر کنترل نماید.

کلمات کلیدی: آونگ معکوس دورانی، کنترل هوشمند، یادگیری تقویتی، یادگیری تقویتی گسسته، یادگیری تقویتی پیوسته

فهرست مطالب

۹	۱- مقدمه.....
۹	۱-۱- تعریف مسئله.....
۱۰	۱-۲- ضرورت اجرای پروژه.....
۱۱	۱-۳- پیشینه‌ی پژوهش.....
۱۱	۱-۳-۱- پیشینه‌ی پژوهش در دنیا.....
۱۲	۱-۳-۲- پیشینه پژوهش در دانشکده مکانیک دانشگاه تهران.....
۱۳	۲- معادلات ریاضی حاکم بر سیستم.....
۱۳	۲-۱- دینامیک آونگ معکوس.....
۱۷	۲-۲- موتور DC.....
۱۸	۳- اصلاحات سازه‌ی آونگ معکوس دورانی.....
۱۸	۳-۱- طراحی سیستم کالیراسیون برای انکودرها.....
۲۱	۳-۲- تغییر میکروکنترلرها.....
۲۲	۳-۳- طراحی برد مدار چاپی.....
۲۲	۳-۳-۱- برد نصب‌شده روی پایه.....
۲۴	۳-۳-۲- برد نصب‌شده روی بازو.....
۲۶	۳-۴- تغییر شیوه‌ی اتصال بازو به پایه.....
۲۷	۳-۵- طراحی مکانیزم رگلاژ برای تسمه تایم متصل به موتور.....
۲۸	۳-۶- لرزش سازه در هنگام کار موتور در سرعت‌های بالا.....
۳۱	۴- محیط شبیه‌سازی.....
۳۱	۴-۱- مدل‌سازی در زبان پایتون.....
۳۲	۴-۱-۱- مفروضات مدل‌سازی.....
۳۲	۴-۱-۲- ورودی و خروجی مدل.....
۳۳	۴-۱-۳- پارامترهای مدل‌سازی.....
۳۳	۴-۲- مدل‌سازی در سیمولینک.....
۳۴	۴-۲-۱- مفروضات مدل‌سازی.....
۳۴	۴-۲-۲- ورودی و خروجی‌های مدل.....
۳۴	۴-۲-۳- پارامترهای مدل‌سازی.....
۳۵	۴-۲-۴- به‌دست آوردن ضریب اصطکاک ویسکوز شفت.....
۳۷	۴-۲-۵- به‌دست آوردن پارامترهای موتور (مکانیزم اندازه‌گیری مکان آن).....
۴۲	۴-۲-۶- به‌دست آوردن مشخصه انتقالی درایور موتور.....
۴۴	۴-۲-۷- پیاده‌سازی سیستم در سیمولینک.....
۴۶	۴-۲-۸- کار با محیط شبیه‌سازی.....

۴۷.....	۵- یادگیری تقویتی.....
۴۷.....	۵-۱- تعاریف پایه.....
۴۹.....	۵-۲- هدف نهایی مسائل یادگیری تقویتی.....
۵۰.....	۵-۳- مسائل یادگیری تقویتی گسسته.....
۵۰.....	۵-۳-۱- فرمولاسیون ریاضی مسائل یادگیری گسسته.....
۵۰.....	۵-۳-۱-۱- فرایند تصمیم‌گیری مارکوف.....
۵۱.....	۵-۳-۱-۲- بهینه‌سازی در تصمیم‌گیری مارکوف.....
۵۱.....	۵-۳-۱-۳- مفهوم State Value و Action Value.....
۵۲.....	۵-۳-۱-۴- معادلات بلمن.....
۵۳.....	۵-۳-۲- شیوه‌ی پداده‌سازی تابع سیاست در محیط‌های گسسته.....
۵۳.....	۵-۳-۳- یادگیری در محیط‌های گسسته.....
۵۵.....	۵-۳-۳-۱- الگوریتم Q-Learning.....
۵۶.....	۵-۳-۳-۲- الگوریتم Deep Q-Learning.....
۵۷.....	۵-۴- مسائل یادگیری تقویتی پیوسته.....
۵۸.....	۵-۴-۱- فرایند یادگیری در محیط پیوسته.....
۵۸.....	۵-۴-۱-۱- الگوریتم Deep Deterministic Policy Gradient (DDPG).....
۶۱.....	۵-۵- کاربرد یادگیری تقویتی در مسائل حوزه مهندسی کنترل.....
۶۲.....	۶- کنترل آونگ معکوس دورانی به کمک یادگیری تقویتی.....
۶۲.....	۶-۱- ورودی الگوریتم کنترلی.....
۶۳.....	۶-۲- خروجی الگوریتم کنترلی.....
۶۴.....	۷- فرایند آموزش.....
۶۵.....	۷-۱- الگوریتم Deep Q-Learning.....
۶۶.....	۷-۲- الگوریتم Deep Deterministic Policy Gradient.....
۶۷.....	۸- نتایج آموزش.....
۶۹.....	۹- تفسیر و بررسی نتایج.....
۶۹.....	۱۰- پیشنهادهایی جهت بهبود نتایج.....
۷۱.....	۱۱- جمع‌بندی و نتیجه‌گیری.....
۷۲.....	مراجع.....
۷۳.....	پیوست‌ها.....
۷۴.....	پیوست ۱: نقشه کارگاهی قطعات تغییر یافته در آونگ معکوس دورانی.....
۸۱.....	پیوست ۲: تصاویر PCBهای طراحی شده در محیط Altium Designer.....

فهرست اشکال

- شکل ۱: شماتیک آونگ معکوس دورانی ۹
- شکل ۲: Quanser Qube برای کنترل آونگ معکوس دورانی ۱۱
- شکل ۳: نمونه اولیه آونگ معکوس دورانی ساخته شده ۱۲
- شکل ۴: آونگ معکوس دورانی ۱۳
- شکل ۵: مدل موتور DC آرمیچر ۱۷
- شکل ۶: دیسک دوار متصل به انکودر ۱۹
- شکل ۷: گیرنده و فرستنده مادون قرمز مدل GK152 ۱۹
- شکل ۸: مدار تبدیل خروجی اپتوکانتیر به مقادیر باینری ۲۰
- شکل ۹: مکانیزم کالیبراسیون انکودر متصل به آونگ ۲۰
- شکل ۱۰: ماژول Womos D1 دارای هسته ESP8266 جهت استفاده روی پایه سازه (Base ESP) ۲۱
- شکل ۱۱: ماژول Womos D1 دارای هسته ESP8266 و جا باتری و شارژر باتری لیتیوم-یون ۱۸۶۵۰ جهت استفاده روی بازوی سازه (Arm ESP) ۲۱
- شکل ۱۲: ارتباط بین بوردهای ESP و سرور محاسباتی ۲۲
- شکل ۱۳: مدار کاهنده ولتاژ برای PCB نصب شده روی پایه سازه ۲۳
- شکل ۱۴: شماتیک مدار چاپی نصب شده بر روی پایه سازه ۲۳
- شکل ۱۵: تصویر سه بعدی مدار چاپی نصب شده روی پایه در نرم افزار Altium Designer ۲۴
- شکل ۱۶: مدار چاپی متصل به پایه جدا از سیستم ۲۴
- شکل ۱۷: مدار چاپی متصل به پایه متصل ۲۴
- شکل ۱۸: شماتیک مدار چاپی نصب شده روی بازوی سازه ۲۵
- شکل ۱۹: تصویر سه بعدی مدار چاپی نصب شده روی بازو در نرم افزار Altium Designer ۲۵
- شکل ۲۰: پولی متصل به شفت موتور و بازوی سازه ی آونگ معکوس دورانی ۲۶
- شکل ۲۱: اتصال بازو به موتور در سازه ی آونگ دورانی معکوس ۲۶
- شکل ۲۲: نمایی از مکانیزم رگلاژ تسمه تایم ۲۷
- شکل ۲۳: مکانیزم تغییر مکان انکودر جهت رگلاژ تسمه تایم ۲۷
- شکل ۲۴: تسمه تایم متصل کننده شفت های موتور و انکودر حالت آزاد ۲۸
- شکل ۲۵: تسمه تایم متصل کننده شفت های موتور و انکودر حالت سفت ۲۸
- شکل ۲۶: سازه ی آونگ معکوس دورانی - نمای روبرو ۲۹
- شکل ۲۷: سازه ی آونگ معکوس دورانی - نمای روبرو ۳۰
- شکل ۲۸: پنجره گرافیکی نوشته شده به زبان پایتون جهت نمایش وضعیت لحظه ای آونگ معکوس دورانی شبیه سازی شده ۳۲
- شکل ۲۹: پنجره گرافیکی سمولینک جهت نمایش وضعیت لحظه ای آونگ معکوس دورانی شبیه سازی شده ۳۳
- شکل ۳۰: آونگ ساده ۳۵

شکل ۳۱: نمودار مکان زاویه‌ای برحسب زمان برای آونگ در هنگام رها شدن از زاویه θ_0	۳۶
شکل ۳۲: برازش منحنی نمودار مکان زاویه‌ای برحسب زمان برای آونگ در هنگام رها شدن از زاویه θ_0	۳۷
شکل ۳۳: نمودار سرعت دوران شفت موتور در هنگام چرخاندن آن به وسیله دریل.....	۳۸
شکل ۳۴: نمودار ولتاژ اندازه‌گیری شده در پایه‌های موتور در هنگام چرخاندن شفت آن به وسیله دریل.....	۳۹
شکل ۳۵: مقادیر به دست آمده برای Kb در طول زمان.....	۳۹
شکل ۳۶: پاسخ پله‌ی موتور به ازای ولتاژ 24V.....	۴۰
شکل ۳۷: پاسخ پله‌ی موتور در کنار پاسخ پله‌ی تابع تبدیل شناسایی شده برای آن.....	۴۱
شکل ۳۸: شمای کلی سیستم آونگ معکوس دورانی در سیمولینک.....	۴۴
شکل 39: نمایی از زیرسیستم <code>rotary_inverted_pendulum_body</code>	۴۵
شکل ۴۰: نمایی از زیرسیستم <code>dc_motor</code>	۴۵
شکل ۴۱: نمایی از زیرسیستم <code>motor_driver</code>	۴۶
شکل ۴۲: شماتیک فرایندهای یادگیری تقویتی.....	۴۷
شکل ۴۳: یک فرایند تصمیم‌گیری مارکوف ساده با سه حالت (State) و دو عمل (Action).....	۵۰
شکل ۴۴: فرایند تکراری ارزیابی و بهبود سیاست.....	۵۴
شکل ۴۵: اجرای هم‌زمان ارزیابی و بهبود سیاست.....	۵۴
شکل ۴۶: تشابه یادگیری تقویتی و فرایند کنترل با بازخورد.....	۶۱
شکل ۴۷: نتایج آموزش الگوریتم Deep Q-Learning در ۴۰۰ اپیزود پایانی.....	۶۷
شکل ۴۸: نتایج آموزش الگوریتم Deep Deterministic Policy Gradient در ۴۰۰ اپیزود پایانی.....	۶۷
شکل ۴۹: نتایج آموزش الگوریتم Deep Q-Learning در ۱۰۰ اپیزود آموزش روس سازی آونگ معکوس دورانی.....	۶۸
شکل ۵۰: نتایج آموزش الگوریتم Deep Deterministic Policy Gradient در ۱۰۰ اپیزود آموزش روس سازی آونگ معکوس دورانی.....	۶۸

فهرست جداول

جدول ۱: پارامترهای دخیل در دینامیک آونگ معکوس دورانی.....	۱۳
جدول ۲: پارامترهای دخیل در مدل‌سازی موتور DC.....	۱۷
جدول ۳: ورودی‌ها و خروجی‌های مدل اولیه.....	۳۲
جدول ۴: پارامترهای دخیل در مدل‌سازی اولیه.....	۳۳
جدول ۵: پارامترهای دخیل در مدل‌سازی ثانویه.....	۳۴
جدول ۶: مقادیر پارامترهای دخیل در مدل‌سازی بازو.....	۴۲
جدول ۷: مقادیر پارامترهای دخیل در مدل‌سازی آونگ.....	۴۳

جدول ۸: مقادیر پارامترهای دخیل در مدل‌سازی موتور DC (و انکودر متصل به آن).....	۴۳
جدول ۹: مقادیر پارامترهای دخیل در مدل‌سازی انکودر متصل به آونگ.....	۴۳
جدول ۱۰: مقادیر پارامترهای دخیل در مدل‌سازی مورد چاپی متصل به بازو.....	۴۳
جدول ۱۱: نمادگذاری متداول برای مفاهیم پایه‌ای در یادگیری تقویتی.....	۴۹
جدول ۱۲: تشابه مفاهیم حوزه‌ی کنترل و یادگیری تقویتی.....	۶۱

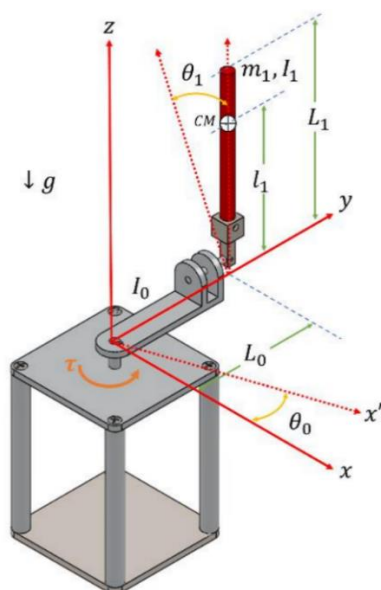
فهرست شبه‌کدها

شبه کد ۱: الگوریتم Q-Learning.....	۵۵
شبه کد ۲: الگوریتم Deep Q-Learning.....	۵۷
شبه کد ۳: الگوریتم Deep Deterministic Policy Gradient (DDPG).....	۶۰

۱- مقدمه

۱-۱- تعریف مسئله

آونگ معکوس دورانی^۱ یک سیستم مکانیکی با دو درجه آزادی و یک محرک است. این سیستم، شامل یک بازوی دوار^۲ است که به وسیله یک موتور حرکت دورانی می‌کند. بر لبه بازو، یک آونگ^۳ متصل شده است که در صفحه‌ی عمود بر بازو حرکت دورانی می‌کند. شکل ۱، طرح شماتیک آونگ معکوس دورانی را نشان می‌دهد:



شکل ۱: شماتیک آونگ معکوس دورانی [1]

هدف از کنترل آونگ معکوس دورانی، نگهداشتن آونگ در وضعیت قائم (به صورت معکوس) به کمک حرکت دورانی موتور است. دینامیک غیرخطی و ناپایدار آونگ معکوس دورانی موجب اهمیت آن به عنوان یکی از مسائل معیار^۴ در حوزه کنترل شده است. در سال‌های اخیر، کنترل پایداری مکانیزم‌های مبتنی بر آونگ معکوس کاربرد گسترده‌ای در صنایع مختلف یافته است. کنترل پایداری موتورسیکلت‌های بدون سرنشین، تخته‌روها (مانند Segway) و ربات‌های دوچرخ و کروی از نمونه‌های کاربرد این فناوری به شمار می‌آیند.

¹ Rotary Inverted Pendulum

² Arm

³ Pendulum

⁴ Benchmark

۱-۲- ضرورت اجرای پروژه

در سال‌های اخیر، در پی پیشرفت‌های چشمگیر در هوش مصنوعی و یادگیری ماشینی، روش‌های هوشمند در حوزه‌ی کنترل محبوبیت فراوانی یافته‌اند. توانایی این روش‌ها در فراگیری روش کنترل محیط‌هایی پیچیده، بدون نیاز به آگاهی از معادلات و دینامیک آن‌ها، همواره موردتوجه علاقه‌مندان این حوزه بوده است. بدون شک یکی از مهم‌ترین ویژگی الگوریتم‌های هوشمند، توانایی هماهنگ شدنشان با تغییرات ناخواسته در سیستم و مقاومت بالایشان در مقابل عدم قطعیت‌های موجود است. در مقابل، بزرگ‌ترین ایراد وارده به این روش‌ها، پیچیدگی پیاده‌سازی آن‌ها و فرایند آموزش طولانی می‌باشد. به همین علت، قبل از جایگزینی روش‌های کنترل کلاسیک، الگوریتم‌های هوشمند راه طولانی در پیش دارند. تحقیقات بسیاری در زمینه‌ی کاهش پیچیدگی و مدت‌زمان یادگیری الگوریتم‌های هوشمند انجام شده است و بسیاری از موفقیت‌های این زمینه، ضمن به چالش کشیدن الگوریتم‌های هوشمند روی مسائل خاص صورت گرفته است.

مسئله‌ی کنترل آونگ معکوس دورانی (Rotary Inverted Pendulum)، به علت دینامیک پیچیده و به‌شدت ناپایداری، همواره موردتوجه پژوهشگران حوزه‌ی کنترل قرار گرفته تا آنجا که از آن به‌عنوان یک مسئله‌ی معیار جهت مقایسه‌ی روش‌های مختلف کنترل یاد می‌شود. از این تئوری موجود در این سیستم، برای مدل‌سازی وسایل گوناگونی از جمله سگوی^۱ و ربات خود تعادل^۲ استفاده می‌شود. در این پروژه، قصد داریم عملکرد کنترلرهای مبتنی بر یادگیری تقویتی^۳ پیوسته و گسسته را بر روی سیستم آونگ معکوس دورانی بسنجیم. به‌عنوان یک مسئله‌ی پیچیده و معیار، انتظار می‌رود که نتایج حاصل از این پژوهش دبدی دقیق و کامل از نقاط ضعف و قوت الگوریتم‌های یادگیری تقویتی به ما بدهد.

گام بعدی این پژوهش، می‌تواند استفاده از نتایج به‌دست‌آمده جهت بهبود عملکرد این الگوریتم‌ها در مسئله‌ی کنترل آونگ معکوس دورانی باشد. این کار نه‌تنها منجر به طراحی یک سیستم کنترلی بهتر مبتنی بر یادگیری تقویتی برای آونگ معکوس دورانی شود، می‌تواند به‌طور غیرمستقیم به توسعه‌ی روش‌های طراحی سیستم‌های یادگیری تقویتی نیز کمک شایانی نماید.

¹ Segway

² Self-Balancing Robot

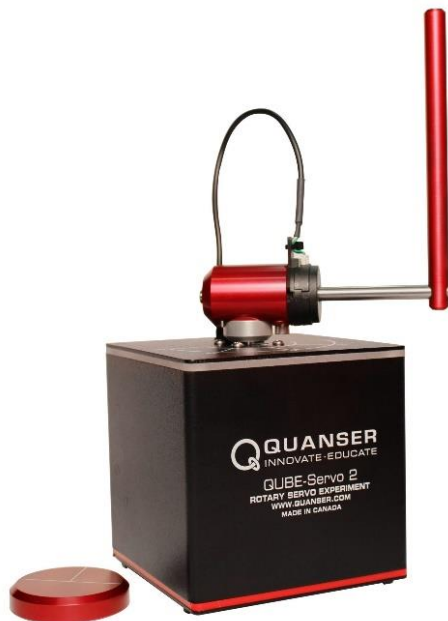
³ Reinforcement Learning

۳-۱- پیشینه‌ی پژوهش

۱-۳-۱- پیشینه‌ی پژوهش در دنیا

مقاله‌های گوناگونی پیرامون کنترل سیستم آونگ معکوس با روش‌های مختلف کنترلی در مجله‌های علمی به چاپ رسیده است. از میان این روش‌ها می‌توان به کنترلرهای خانواده PID [2] [3]، کنترل فازی [4]، روش LQR [5] و یا روش‌های هوشمند مانند الگوریتم ژنتیک، الگوریتم ازدحام ذرات^۱ [6] و ... اشاره نمود. باوجود بررسی عملکرد تعداد زیادی از روش‌های کنترلی هوشمند جهت کنترل آونگ معکوس دورانی، اکثر پژوهش‌های انجام شده بر روی مدل‌سازی کامپیوتری تمرکز دارند و در کمتر مقاله‌ای به بررسی عملکرد کنترلر بر روی سیستم فیزیکی پرداخته شده است.

درزمینه‌ی ساخت و کنترل مکانیزم آونگ معکوس، شرکت کوآنسر^۲، از شرکت‌های مطرح درزمینه‌ی ساخت لوازم آزمایشگاهی، فعال می‌باشد. این شرکت محصولی به نام Quanser Qube تولید کرده که قابلیت تبدیل شدن به تعدادی از مسائل کلاسیک و معیار حوزه‌ی کنترل، ازجمله مسئله‌ی آونگ معکوس دورانی، را دارد.



شکل ۲: Quanser Qube برای کنترل آونگ معکوس دورانی^۳

^۱ Particle Swarm Optimization (PSO)

^۲ Quanser

^۳ برگرفته از Quanser.com

۱-۳-۲- پیشینه پژوهش در دانشکده مکانیک دانشگاه تهران

پروژه ساخت آونگ معکوس دورانی از اواخر سال ۱۳۹۸ در آزمایشگاه کاربردهای هوش مصنوعی در مکانیک (AIME¹ Lab.) به هدف بررسی عملکرد انواع روش‌های کنترلی بر روی یک مسئله معیار^۲ در حوزه کنترل آغاز شد. در گام نخست، سازه‌ی آونگ معکوس دورانی ساخته شد و جهت بررسی عملکرد کنترلر PID مورد مطالعه قرار گرفت [7]. تصویر زیر، نمایی از آونگ دورانی معکوسی ساخته شده می‌باشد:



شکل ۳: نمونه اولیه آونگ معکوس دورانی ساخته شده

در پژوهش بعدی به بررسی عملکرد الگوریتم Q-Learning گسسته (از الگوریتم‌های حوزه یادگیری تقویتی^۳) جهت پایدار کردن آونگ معکوس دورانی پرداخته شد و نتایج به دست آمده از الگوریتم هوشمند با نتایج پیشین به دست آمده از کنترل PID مقایسه گردید [8]. متأسفانه، برخی از مشکلات موجود در سازه باعث ایجاد اختلال در روند آموزش آونگ شد و عملکرد نهایی الگوریتم Q-Learning از مقدار مورد انتظار ضعیف تر بود.

¹ Artificial Intelligence in Mechanical Engineering

² Benchmark

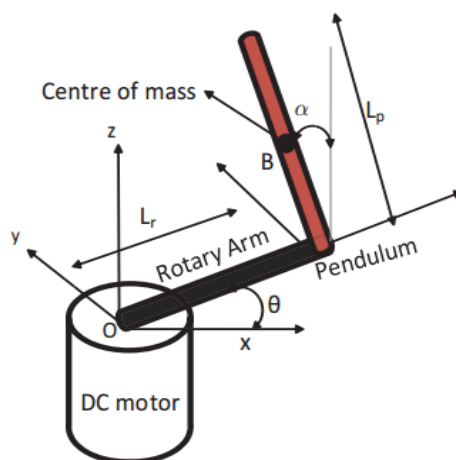
³ Reinforcement Learning

۲- معادلات ریاضی حاکم بر سیستم

در این بخش به بررسی قوانین و معادلات فیزیکی حاکم بر سازه‌ی آونگ معکوس دورانی خواهیم پرداخت. اگرچه استفاده از این روابط در کنترل سازه به کمک الگوریتم‌های یادگیری تقویتی ضروری نیست، دانستن این معادلات برای فرایند شبیه‌سازی^۱ بسیار مهم و حائز اهمیت می‌باشد.

۱-۲- دینامیک آونگ معکوس

شکل زیر، شماتیک آونگ معکوس دورانی را نشان می‌دهد:



شکل ۴: آونگ معکوس دورانی^۲

همچنین، پارامترهای دخیل در مدل‌سازی دینامیک آونگ به شرح موجود در جدول زیر می‌باشند:

جدول ۱: پارامترهای دخیل در دینامیک آونگ معکوس دورانی

واحد	معنی	نماد
rad	مکان زاویه‌ای بازو	θ
rad	مکان زاویه‌ای آونگ	α
m	مؤلفه‌های مکان انتهای بازو	$[x_r, y_r, z_r]$
m	مؤلفه‌های مکان مرکز جرم آونگ	$[x_p, y_p, z_p]$
kg	جرم آونگ	m_p
$kg.m^2$	اینرسی آونگ	J_p
m	طول آونگ	L_p
$N.m.s/rad$	ضریب میرایی ویسکوز اتصال آونگ و بازو	c_p
m	طول بازوی دوار	L_r
$kg.m^2$	اینرسی بازوی دوار	J_a
$N.m.s/rad$	ضریب میرایی ویسکوز اتصال بازو و موتور	c_a

¹ Simulation

^۲ برگرفته از degruyter.com

در ادامه، معادلات دینامیکی حاکم بر آونگ را به دست می‌آوریم.

برای نقطه‌ی محل اتصال آونگ به بازو داریم:

$$\begin{bmatrix} x_r \\ y_r \end{bmatrix} = \begin{bmatrix} L_r \cos \theta \\ L_r \sin \theta \end{bmatrix} \quad (\text{Equation 1})$$

با مشتق‌گیری از معادله‌ی ۱ سرعت انتهای بازو به‌صورت زیر به دست می‌آید.

$$\begin{bmatrix} \dot{x}_r \\ \dot{y}_r \end{bmatrix} = \begin{bmatrix} -L_r \dot{\theta} \sin \theta \\ L_r \dot{\theta} \cos \theta \end{bmatrix} \quad (\text{Equation 2})$$

با فرض توزیع یکنواخت جرم آونگ در طول آن، مرکز جرم آونگ در وسط آن قرار می‌گیرد. برای این نقطه می‌توان نوشت:

$$\begin{bmatrix} x_p \\ y_p \\ z_p \end{bmatrix} = \begin{bmatrix} L_r \cos \theta + \frac{L_p \sin \alpha \sin \theta}{2} \\ L_r \sin \theta - \frac{L_p \sin \alpha \cos \theta}{2} \\ \frac{L_p}{2} \cos \alpha \end{bmatrix} \quad (\text{Equation 3})$$

با مشتق‌گیری از معادله‌ی ۳ سرعت مرکز جرم آونگ به‌صورت زیر به دست می‌آید.

$$\begin{bmatrix} \dot{x}_p \\ \dot{y}_p \\ \dot{z}_p \end{bmatrix} = \begin{bmatrix} -L_r \dot{\theta} \sin \theta + \frac{L_p \dot{\alpha} \cos \alpha \sin \theta}{2} + \frac{L_p \dot{\theta} \sin \alpha \cos \theta}{2} \\ L_r \dot{\theta} \sin \theta - \frac{L_p \dot{\alpha} \cos \alpha \cos \theta}{2} + \frac{L_p \dot{\theta} \sin \alpha \sin \theta}{2} \\ \frac{L_p}{2} \dot{\alpha} \sin \alpha \end{bmatrix} \quad (\text{Equation 4})$$

نیروی جنبشی کل سیستم از جمع انرژی جنبشی آونگ و بازوی دوار به دست می‌آید که به‌صورت زیر است:

$$KE = KE_{arm} + KE_{pend} \quad (\text{Equation 5})$$

با استفاده از معادله‌ی ۵ و مقادیر به‌دست‌آمده برای سرعت‌ها داریم:

$$KE = \frac{1}{2} \left(m_p L_r^2 + J_r + m_p \frac{L_p^2}{4} \sin^2 \alpha \right) \dot{\theta}^2 + \frac{1}{2} \left(m_p \frac{L_p^2}{4} + J_p \right) \dot{\alpha}^2 - \frac{1}{2} (m_p L_r L_p \dot{\theta} \dot{\alpha} \cos \alpha) \quad (\text{Equation 6})$$

به‌طور مشابه، برای انرژی پتانسیل کل سیستم داریم:

$$PE = PE_{arm} + PE_{pend} \quad (\text{Equation 7})$$

از آنجاکه بازوی دوار حرکت عمودی ندارد و تغییرات انرژی پتانسیل صفر است، مجموع انرژی جنبشی از رابطه‌ی زیر به دست می‌آید:

$$PE = 0 + \left(-m_p g \left(\frac{L_p}{2} - \frac{L_p}{2} \cos \alpha \right) \right) \quad (\text{Equation 8})$$

پارامتر لاگرانژین^۱ برای سیستم، عبارت است از تفاضل انرژی جنبشی و انرژی پتانسیل سیستم که به‌صورت زیر تعریف می‌شود:

$$\mathcal{L} = KE - PE \quad (\text{Equation 9})$$

با قرار دادن مقادیر به‌دست‌آمده برای انرژی جنبشی و پتانسیل سیستم در معادله‌ی ۹ داریم:

$$\mathcal{L} = \frac{1}{2} \left(m_p L_r^2 + J_r + m_p \frac{L_p^2}{4} \sin^2 \alpha \right) \dot{\theta}^2 + \frac{1}{2} \left(m_p \frac{L_p^2}{4} + J_p \right) \dot{\alpha}^2 - \frac{1}{2} (m_p L_r L_p \dot{\theta} \dot{\alpha} \cos \alpha) + m_p g \left(\frac{L_p}{2} - \frac{L_p}{2} \cos \alpha \right) \quad (\text{Equation 10})$$

تابع رایلی^۲ برای این سیستم عبارت است از:

$$R = \frac{1}{2} c_r \dot{\theta}^2 + \frac{1}{2} c_p \dot{\alpha}^2 \quad (\text{Equation 11})$$

¹ Lagrangian

² Riley Function

به دلیل این که سیستم را می توان به ۲ درجه آزادی تقسیم کرد، در نتیجه ۲ معادله لاگرانژ برای این سیستم وجود دارد که به صورت زیر نوشته می شوند:

$$\frac{d}{dt}\left(\frac{\partial \mathcal{L}}{\partial \dot{\theta}}\right) - \frac{\partial \mathcal{L}}{\partial \theta} + \frac{\partial R}{\partial \dot{\theta}} = \tau \quad , \quad \frac{d}{dt}\left(\frac{\partial \mathcal{L}}{\partial \dot{\alpha}}\right) - \frac{\partial \mathcal{L}}{\partial \alpha} + \frac{\partial R}{\partial \dot{\alpha}} = 0 \quad (\text{Equations 12\&13})$$

حال با قرار دادن پارامترهای به دست آمده از معادلات قبل و خطی سازی سیستم پیرامون نقطه‌ی کاری سیستم ($\alpha = 0$) به کمک سری تیلور^۱ حول صفر، معادلات حرکت سیستم حول نقطه‌ی $\alpha = 0$ به صورت زیر به دست می آید:

$$(m_p L_r^2 + J_r)\ddot{\theta} - \frac{1}{2}m_p L_p L_r \ddot{\alpha} + c_r \dot{\theta} = \tau \quad (\text{Equation 14})$$

$$-\frac{1}{2}m_p L_p L_r \ddot{\theta} + (J_p + \frac{1}{4}m_p L_p^2)\ddot{\alpha} - \frac{m_p L_p g \alpha}{2} + c_p \dot{\alpha} = 0 \quad (\text{Equation 15})$$

با ساده سازی معادلات بالا، روابط زیر برای شتاب زاویه‌ای بازوی دوار و آونگ به دست می آید:

$$\ddot{\theta} = \frac{1}{J_T} \left(- \left(J_p + \frac{1}{4}m_p L_p^2 \right) c_r \dot{\theta} - \frac{1}{2}m_p L_p L_r c_p \dot{\alpha} + \frac{1}{4}m_p^2 L_p^2 L_r g \alpha + \left(J_p + \frac{1}{4}m_p L_p^2 \right) \tau \right) \quad (\text{Equation 16})$$

$$\ddot{\alpha} = \frac{1}{J_T} \left(\frac{1}{2}m_p L_p L_r c_r \dot{\theta} - (J_r + m_p L_r^2) c_p \dot{\alpha} + \frac{1}{2}m_p L_p g \alpha - \frac{1}{2}m_p L_p L_r \tau \right) \quad (\text{Equation 17})$$

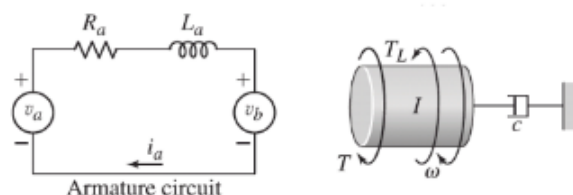
در معادلات ۱۶ و ۱۷، J_T یک پارامتر کمکی است که به صورت زیر تعریف می شود:

$$J_T = J_p m_p L_r^2 + J_r J_p + \frac{1}{4}J_r m_p L_p^2 \quad (\text{Equation 18})$$

¹ Tylor Series

۲-۲- موتور DC

در آونگ معکوس دورانی ساخته شده، از یک موتور DC، به عنوان محرک سیستم استفاده می‌شود. در این بخش، فرمولاسیون ریاضی برای مدل سازی موتور DC آورده شده است. موتور DC را می‌توان به صورت یک سیستم خطی با شماتیک زیر در نظر گرفت:



شکل ۵: مدل موتور DC آرمیچر^۱

همچنین، پارامترهای دخیل در مدل سازی موتور DC به شرح موجود در جدول زیر می‌باشند:

جدول ۲: پارامترهای دخیل در مدل سازی موتور DC

نماد	معنی	واحد
V_a	ولتاژ آرمیچر	<i>Volt</i>
i_a	جریان آرمیچر	<i>Amps</i>
R_a	مقاومت آرمیچر	<i>Ohms</i>
L_a	اندوکتانس آرمیچر	<i>Henry</i>
V_b	ولتاژ Back emf	<i>Volts</i>
T	گشتاور موتور	<i>N.m</i>
I	اینرسی روتور	<i>Kg.m^2</i>
T_L	گشتاور بار	<i>N.m</i>
θ	زاویه شفت موتور	<i>Rad</i>
C	ضریب میرایی	<i>N.m.s/rad</i>
K_t	ضریب گشتاور	<i>N.m/A</i>
K_b	ضریب Back emf	<i>V.s/m</i>

در ادامه، معادلات حاکم بر موتور را می‌نویسیم:

$$R_a i_a + L_a \frac{di_a}{dt} + V_b = V_a \quad (\text{Equation 19}) \quad V_b = K_b \frac{d\theta}{dt} \quad (\text{Equation 20})$$

$$I \frac{d^2\theta}{dt^2} + C \frac{d\theta}{dt} = T - T_L \quad V_b = K_b \frac{d\theta}{dt} \quad (\text{Equation 21}) \quad T = K_t i_a \quad (\text{Equation 22})$$

^۱ برگرفته از ctms.engin.umich.edu

۳- اصلاحات سازه‌ی آونگ معکوس دورانی

فرایند یادگیری در روش‌های کنترل هوشمند، ازجمله یادگیری تقویتی^۱، اکثراً زمان‌بر هستند. به همین جهت، لازم است تا سازه‌ی فیزیکی برای کارکرد چندین ساعته مستحکم و آماده شود. سازه آونگ معکوس پیشین، دارای چندین نقص بود که پیش از هر چیز لازم بود تا اصلاح شوند. برخی از این نقایص در زیر آورده شده‌اند:

- ۱- بازوی سازه‌ی آونگ معکوس به‌وسیله‌ی یک پیچ به‌موازات شفت موتور، به این شفت متصل می‌شد. لذا، جدای از میزان سفت شدن، به علت اینرسی نسبتاً زیاد بازو و آونگ، بعد از کارکرد چنددقیقه‌ای پیچ شل می‌شد.
- ۲- سیستم در هنگام چرخش موتور در بالاترین دور، ارتعاشات شدیدی داشت و ثابت نمی‌ایستاد.
- ۳- انکودر ثابت که به‌وسیله تسمه تایم و پولی به شفت دوار موتور متصل بود، با چسب به بدنه چسبانده شده بود. این موضوع باعث می‌شد امکان رگلاژ تسمه تایم و تنظیم کشش آن وجود نداشته باشد.
- ۴- از آنجاکه انکودرهای استفاده‌شده، از نوع افزایشی^۲ هستند، بعد از مدتی کار کردن از حالت کالیبره خارج می‌شدند و خطای اندازه‌گیری قابل‌توجهی به وجود می‌آمد.
- ۵- عمر باتری لیتیم-یون متصل به مدار متحرک بازو، با توجه به مصرف مدار، کمتر از ۲ ساعت بود و سیستم از پس فرایندهای آموزش طولانی‌مدت برنمی‌آمد.
- ۶- جهت پیاده‌سازی مدارهای الکتریکی از برد برد استفاده‌شده بود و شل بودن اتصالات، امکان بروز انواع خرابی و قطعی در مدار در هنگام کار را تشدید می‌کرد. (علی‌الخصوص در مدار متحرک متصل روی بازو)
- ۷- شفت موتور، به علت کارکرد زیاد، به میزان جزئی از حالت عمودی انحراف داشت.

جهت آماده کردن سازه‌ی موجود برای اجرای الگوریتم‌های کنترل هوشمند، لازم بود تا نواقص فوق تا حد امکان اصلاح شوند. در ادامه، به توضیح مختصر اصلاحات صورت گرفته روی سازه‌ی آونگ می‌پردازیم.

۳-۱- طراحی سیستم کالیبراسیون برای انکودرها

در سازه‌ی پیشین از انکودرهای افزایشی^۲ جهت سنجش موقعیت زاویه‌ای موتور و آونگ در هر لحظه استفاده می‌شد. علت استفاده از این نوع انکودر به جای انکودر مطلق^۴، قیمت بسیار پایین‌تر آن بود. در طرف مقابل، انکودرهای افزایشی، به دلیل نداشتن مبدا مشخص، بعد از مدت زمان طولانی کار کردن از حالت کالیبره خارج می‌شوند و اعداد اشتباهی را نشان می‌دهند. این موضوع علی‌الخصوص در هنگام خواندن زاویه‌ی آونگ بسیار مشکل‌ساز می‌شود زیرا حتی اختلاف اندک زاویه‌ی قرائت شده برای آونگ از حالت مقدار واقعی آن می‌تواند موجب ناپایداری سیستم و سقوط آونگ شود.

¹ Reinforcement Learning

² Incremental Encoder

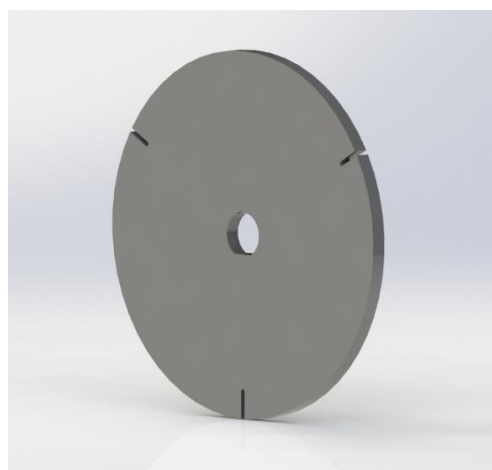
³ Incremental Encoder

⁴ Absolute Encoder

جهت رفع این مشکل یک سیستم کالیبراسیون برای انکودر نصب شده روی بازو (برای اندازه‌گیری زاویه‌ی آونگ) طراحی شد. در این سیستم، از یک اپتوکانترا^۱ مدل GK152 (حاوی یک فرستنده و گیرنده‌ی مادون‌قرمز) و یک دیسک استفاده می‌شود. دیسک دوار به همراه شفت انکودر و آونگ دوران می‌کند و بخشی از آن همواره بین فرستنده و گیرنده‌ی مادون‌قرمز قرار دارد. در سه‌نقطه از مسیر دیسک مقابل حس‌گر، سوراخ‌هایی ایجاد شده است و در هنگام قرار گرفتن این نقاط در بین فرستنده و گیرنده‌ی مادون‌قرمز، نور فرستنده توسط گیرنده دریافت و انکودر کالیبره می‌شود.



شکل ۷: گیرنده و فرستنده مادون‌قرمز مدل GK152



شکل ۶: دیسک دوار متصل به انکودر

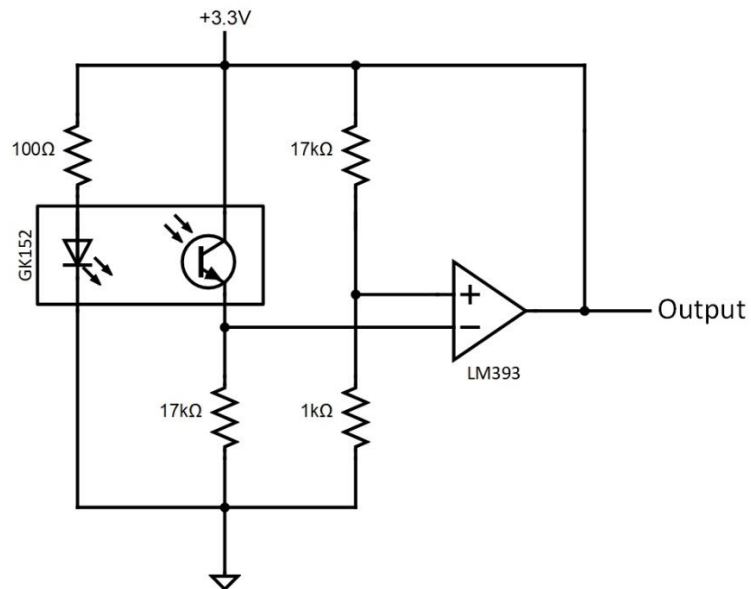
همچنین، از یک تقویت‌کننده‌ی عملیاتی^۲ مقایسه‌گر جهت تبدیل خروجی‌های گیرنده‌ی مادون‌قرمز به مقادیر باینری^۳ استفاده شد. در این مدار، ترانزیستور گیرنده‌ی مادون‌قرمز به‌صورتی در مدار قرار می‌گیرد که بسته به دریافت یا عدم دریافت نور فرستنده، در ناحیه اشباع یا خاموش باشد. درنهایت، ولتاژ پایه امیتر^۴ به کمک یک تقویت‌کننده عملیاتی مقایسه‌گر (LM393) با یک مقدار مرجع مقایسه می‌شود و مقادیر باینری در خروجی تقویت‌کننده عملیاتی تولید می‌شوند. شکل زیر، شماتیک مدار را نشان می‌دهد.

¹ Opto Counter

² Operational Amplifier

³ Binary

⁴ Emitter



شکل ۸: مدار تبدیل خروجی اپتوکانتور به مقادیر باینری

شکل زیر، نمایی از مکانیزم کالیبراسیون انکودر متصل به آونگ را نشان می‌دهد.



شکل ۹: مکانیزم کالیبراسیون انکودر متصل به آونگ

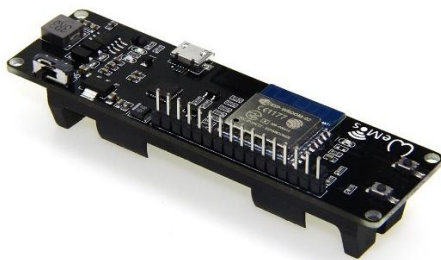
۳-۲- تغییر میکروکنترلرها

در سازه‌ی پیشین از دو برد آردوینو^۱ بر روی بازو و پایه‌ی سازه جهت پیاده‌سازی الگوریتم‌های کنترلی استفاده می‌شد و این دو برد به کمک NRF باهم تبادل اطلاعات می‌کردند. مدار موجود روی بازوی سازه به کمک یک باتری لیتیم و برد نصب شده روی پایه از USB لپ‌تاپ تغذیه می‌شد. ساختار مذکور دو ایراد برجسته داشت:

۱- این ساختار برای پیاده‌سازی کنترلهایی با حجم پردازش پایین مانند کنترلر PID بسیار مناسب است اما توان پردازشی میکروکنترلرهای آردوینو بسیار کمتر از مقدار موردنیاز برای پیاده‌سازی کنترلهای هوشمند مبتنی بر الگوریتم‌های هوش مصنوعی می‌باشد.

۲- NRF و مدار شارژ باتری لیتیم به‌صورت ماژول‌های جداگانه به آردوینو متصل می‌شدند که باعث افزایش تعداد سیم‌کشی در مدار بود.

لذا، آردوینوها با دو ماژول Womos D1 با هسته ESP8266 جایگزین شدند. ماژول نصب شده روی بازوی سازه، به یک مدار شارژ باتری لیتیم مجهز است و نیاز به مدار شارژ مجزا را برطرف می‌سازد.

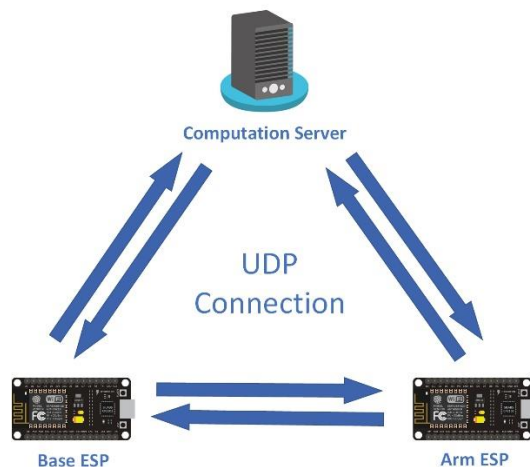


شکل ۱۰: ماژول Womos D1 دارای هسته ESP8266 جهت استفاده روی پایه سازه (Base ESP)
شکل ۱۱: ماژول Womos D1 دارای هسته ESP8266 و جا باتری و شارژر باتری لیتیم-یون ۱۸۶۵۰ جهت استفاده روی بازوی سازه (Arm ESP)

این دو برد، علاوه بر برخورداری از قیمتی پایین به نسبت سایر مدل‌های مشابهشان (ازجمله مدل‌های مجهز به بلوتوث^۲)، به کمک WiFi به یکدیگر و به یک سرور محاسباتی روی کامپیوتر متصل می‌شوند. اتصال به سرور محاسباتی امکان اجرای الگوریتم‌هایی با نیاز به توان محاسباتی بالاتر را مهیا می‌کنند. همچنین، در الگوریتم‌های با توان محاسباتی پایین به‌طور مستقیم روی خود ماژول‌ها قابل اجرا هستند.

^۱ Arduino

^۲ Bluetooth



شکل ۱۲: ارتباط بین بوردهای ESP و سرور محاسباتی

۳-۳- طراحی برد مدار چاپی^۱

مدارات الکتریکی سازه‌ی پیشین، بر روی برد^۲ پیاده‌سازی شده بودند و برطرف کردن مشکلات احتمالی در فرایند آموزش (مانند قطعی‌های احتمالی در مدار) به علت پیچیدگی سیم‌کشی، بسیار سخت و زمان‌بر بود. جهت بهبود عملکرد مدارات الکتریکی در هنگام فرایندهای آموزش طولانی‌مدت، دو برد مدار چاپی مجزا برای سازه طراحی و ساخته شد. در ادامه، به بررسی این مدارها می‌پردازیم:

۳-۳-۱- برد نصب‌شده روی پایه

مدار قرار گرفته روی پایه‌ی سازه، دارای بخش‌های زیر است:

- ۱- پورت اتصال به درایور موتور^۳ BTS7960
- ۲- پورت اتصال Womos D1
- ۳- پورت اتصال انکودر
- ۴- پورت اتصال اپتوکانتور^۴ به‌عنوان مرجع انکودر
- ۵- پورت تغذیه مدار
- ۶- مدار تغذیه Womos D1
- ۷- مدار تبدیل خروجی اپتوکانتور به مقادیر باینری
- ۸- Debugger LED

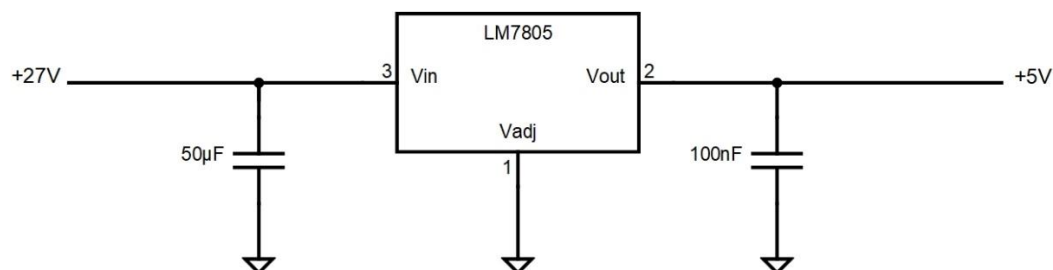
^۱ Printed Circuit Board (PCB)

^۲ Bread Board

^۳ Motor Driver

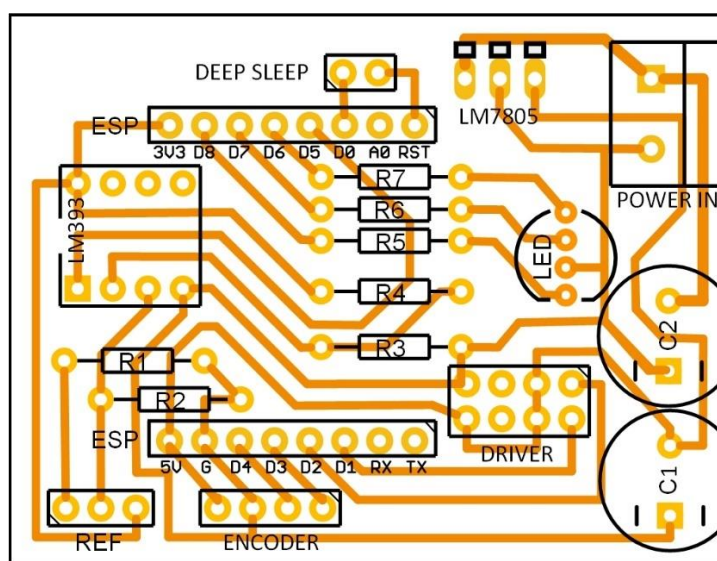
^۴ Opto Counter

این مدار، از برق ۲۷ ولت دی‌سی خروجی پاور سوئیچ^۱ تغذیه می‌شود. سطح ولتاژ ۲۷ ولت به کمک یک رگلاتور ولتاژ^۲ مدل 7805 به ۵ ولت کاهش می‌یابد که برای اجزای مدار قابل استفاده است.



شکل ۱۳: مدار کاهنده ولتاژ برای PCB نصب شده روی پایه سازه

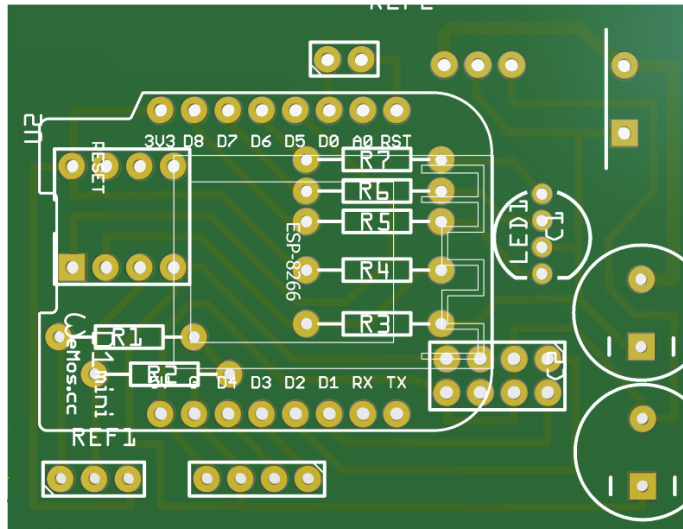
شماتیک مدار چاپی طراحی شده و تصویر سه بعدی آن در نرم افزار Altium Designer در دو شکل بعدی آمده است:



شکل ۱۴: شماتیک مدار چاپی نصب شده بر روی پایه سازه

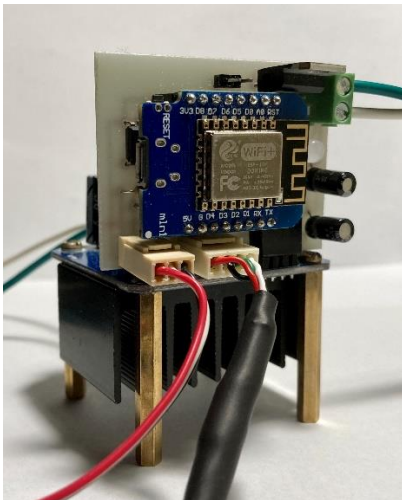
¹ Power Switch

² Voltage Regulator

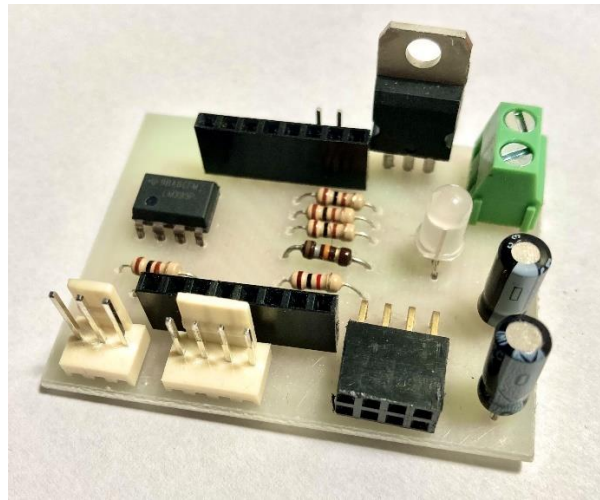


شکل ۱۵: تصویر سه بعدی مدار چاپی نصب شده روی پایه در نرم افزار Altium Designer

همچنین، مدار چاپی نصب شده روی پایه در حالت جدا از سیستم و متصل به آن در ادامه آورده شده است.



شکل ۱۷: مدار چاپی متصل به پایه متصل



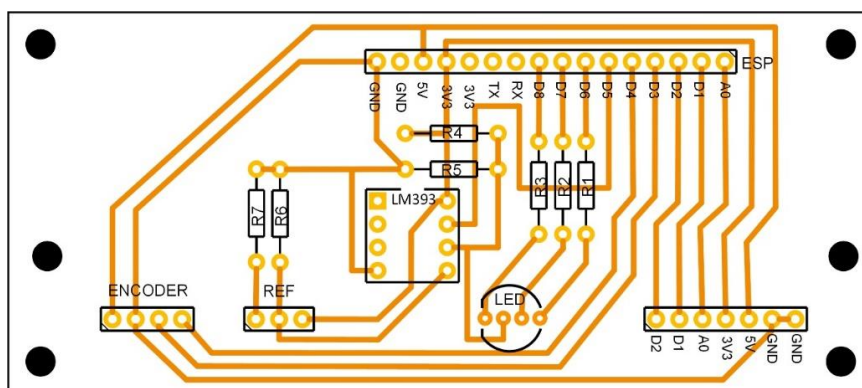
شکل ۱۶: مدار چاپی متصل به پایه جدا از سیستم

۳-۳-۲- مورد نصب شده روی بازو

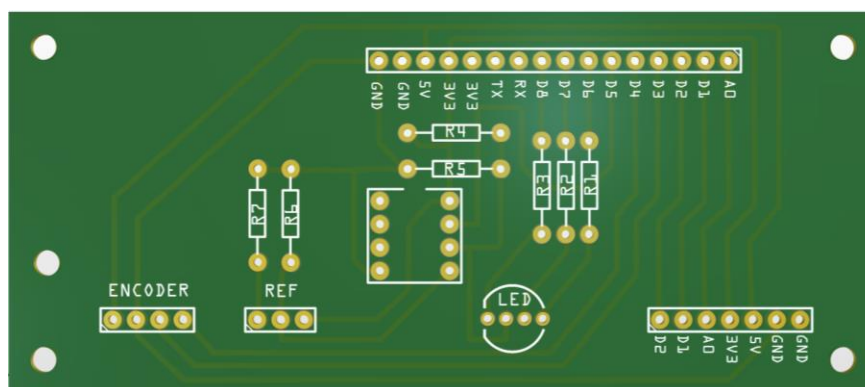
مدار قرار گرفته روی بازوی سازه، دارای بخش‌های زیر است:

- ۱- پورت اتصال Womos D1
- ۲- پورت اتصال انکودر
- ۳- پورت اتصال اپتوکانتربه عنوان مرجع انکودر
- ۴- مدار تبدیل خروجی اپتوکانتربه مقادیر باینری
- ۵- Debugger LED

شماتیک مدار چاپی طراحی شده و تصویر سه بعدی آن در نرم افزار Altium Designer در دو شکل بعدی آمده است:

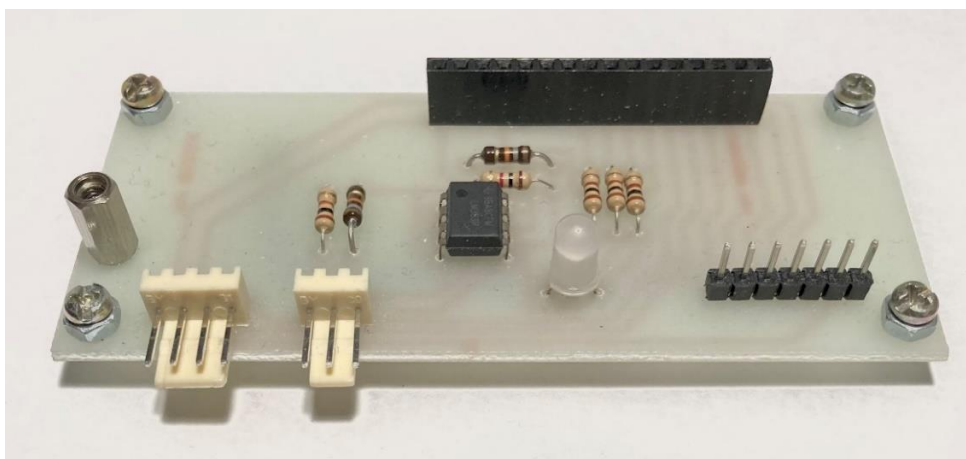


شکل ۱۸: شماتیک مدار چاپی نصب شده روی بازوی سازه

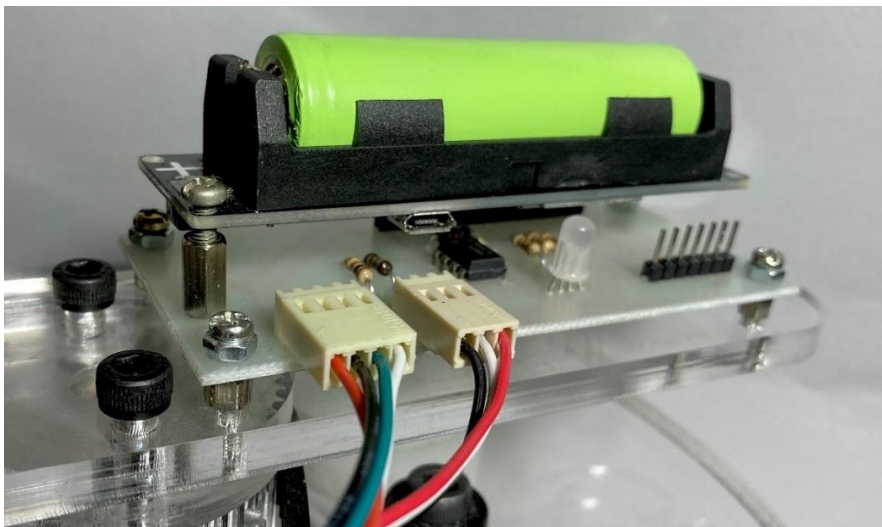


شکل ۱۹: تصویر سه بعدی مدار چاپی نصب شده روی بازو در نرم افزار Altium Designer

همچنین، مدار چاپی نصب شده روی بازو در حالت جدا از سیستم و متصل به آن در ادامه آورده شده است.



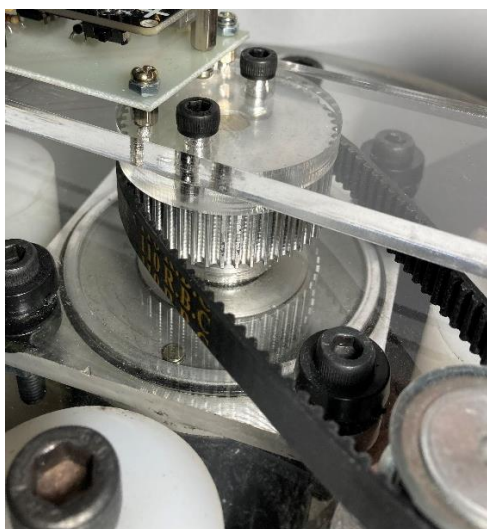
شکل ۱۸: مدار چاپی متصل به بازو در حالت جدای از سیستم



شکل ۱۷: مدار چاپی متصل به بازو در حالت نصب شده روی سیستم

۳-۴- تغییر شیوهی اتصال بازو به پایه

در سازهی آونگ معکوس دورانی پیشین، شفت موتور تنها با استفاده از یک پیچ به موازات آن، به بازوی سازه متصل شده بود. بدون توجه به میزان سفتی پیچ، اتصال بین شفت موتور و بازو بعد از مدتی کارکرد شل می‌شد و لازم بود تا پیچ دوباره سفت شود. برای حل این مشکل، اتصال بین آونگ و موتور با یک اتصال فلنجی جایگزین شد. بدین منظور، یک پولی سری تراش کاری و دو رزوه در طرفین محل قرارگیری شفت موتور در آن ایجاد شد. بر روی بازوی سازه نیز دو سوراخ تعبیه شد و به کمک دو پیچ، بازوی آونگ دورانی معکوس به پولی متصل می‌شود. اتصال پولی به موتور نیز از طریق دو پیچ مغزی صورت می‌گیرد. اشکال زیر، نمایی از پولی تراش کاری شده و شیوهی قرارگیری آن روی سازه است:



شکل ۲۱: اتصال بازو به موتور در سازهی آونگ دورانی معکوس



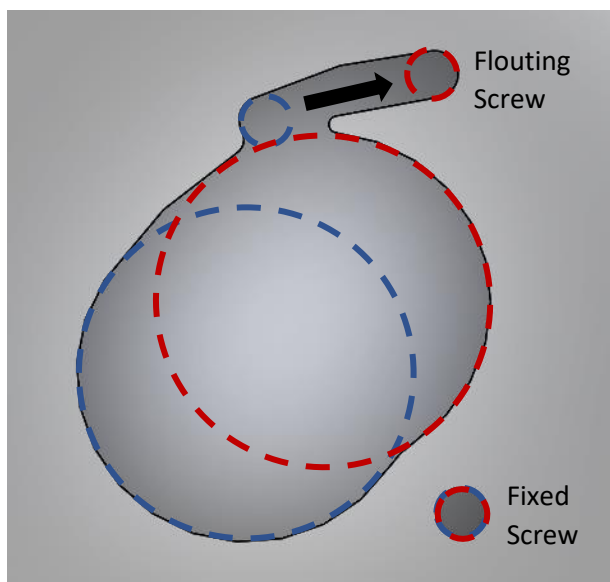
شکل ۲۰: پولی متصل به شفت موتور و بازوی سازهی آونگ معکوس دورانی

نقشه‌ی کارگاهی پولی در ضمیمه ۱ آورده شده است.

۳-۵- طراحی مکانیزم رگلاژ برای تسمه تایم متصل به موتور

جهت اندازه‌گیری مکان زاویه‌ای موتور در هر لحظه، از یک انکودر متصل به پایه سازه استفاده می‌شود. این انکودر به وسیله‌ی تسمه تایم و پولی به موتور متصل می‌شود. در سازه‌ی آونگ معکوس دورانی پیشین، انکودر به وسیله‌ی چسب به پایه‌ی سازه چسبانده شده بود. همچنین، به علت ایجاد امکان تعویض و جایگذاری مجدد تسمه، فاصله‌ی مرکز تا مرکز شفت انکودر تا شفت موتور کمتر از حد استاندارد بود و این موضوع باعث لقی بیش‌ازاندازه‌ی تسمه و ایجاد خطا در اندازه‌گیری می‌شد. به جهت رفع این مشکل، یک مکانیزم رگلاژ ساده برای تغییر فاصله‌ی مرکز شفت‌های انکودر و موتور طراحی شد. این مکانیزم، به انکودر و شفت انکودر به میزان حدوداً ۱ سانتی‌متر اجازه‌ی جابه‌جایی می‌دهد.

تصویر زیر، نمایی از مکانیزم رگلاژ طراحی شده در نرم‌افزار سالی‌دورکس^۱ می‌باشد:



شکل ۲۳: مکانیزم تغییر مکان انکودر جهت رگلاژ تسمه تایم



شکل ۲۲: نمایی از مکانیزم رگلاژ تسمه تایم

نقشه‌ی این مکانیزم رگلاژ تسمه تایم در ضمیمه ۱ آورده شده است.

¹ Solid Works

دو تصویر زیر از تسمه تایم در حالت آزاد و سفت هستند:



شکل ۲۵: تسمه تایم متصل کننده شفت های موتور و انکودر
حالت سفت



شکل ۲۴: تسمه تایم متصل کننده شفت های موتور و انکودر
حالت آزاد

۳-۶- لرزش سازه در هنگام کار موتور در سرعت های بالا

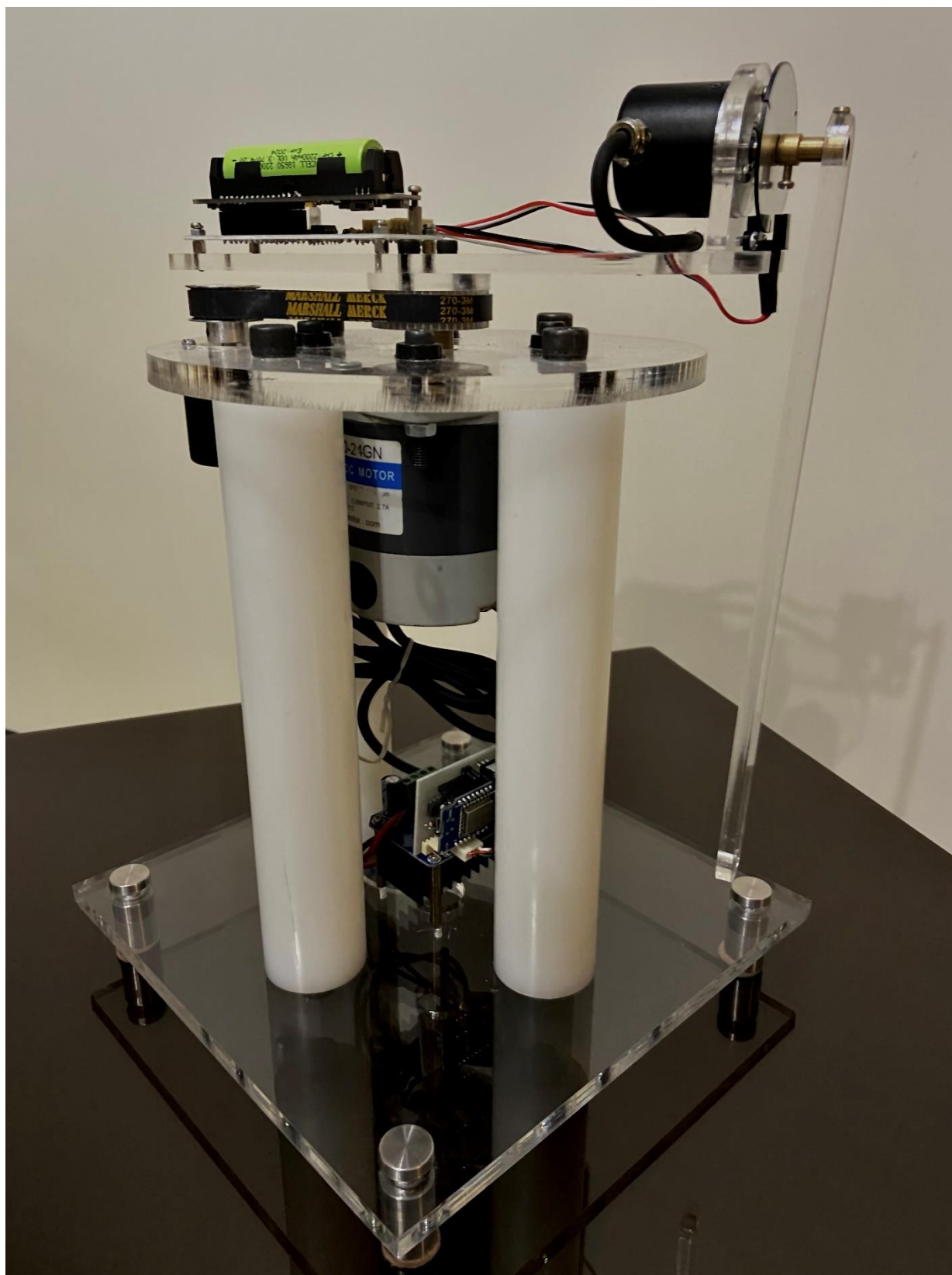
این مشکل عمدتاً ناشی از سه علت بود:

- ۱- وجود نابالانسی در سازه ی بازو و آونگ متصل به موتور
- ۲- عدم وجود پایه مناسب برای دستگاه
- ۳- شیوه ی اتصال موتور به سازه
- ۴- به وجود نابالانسی جزئی در موتور

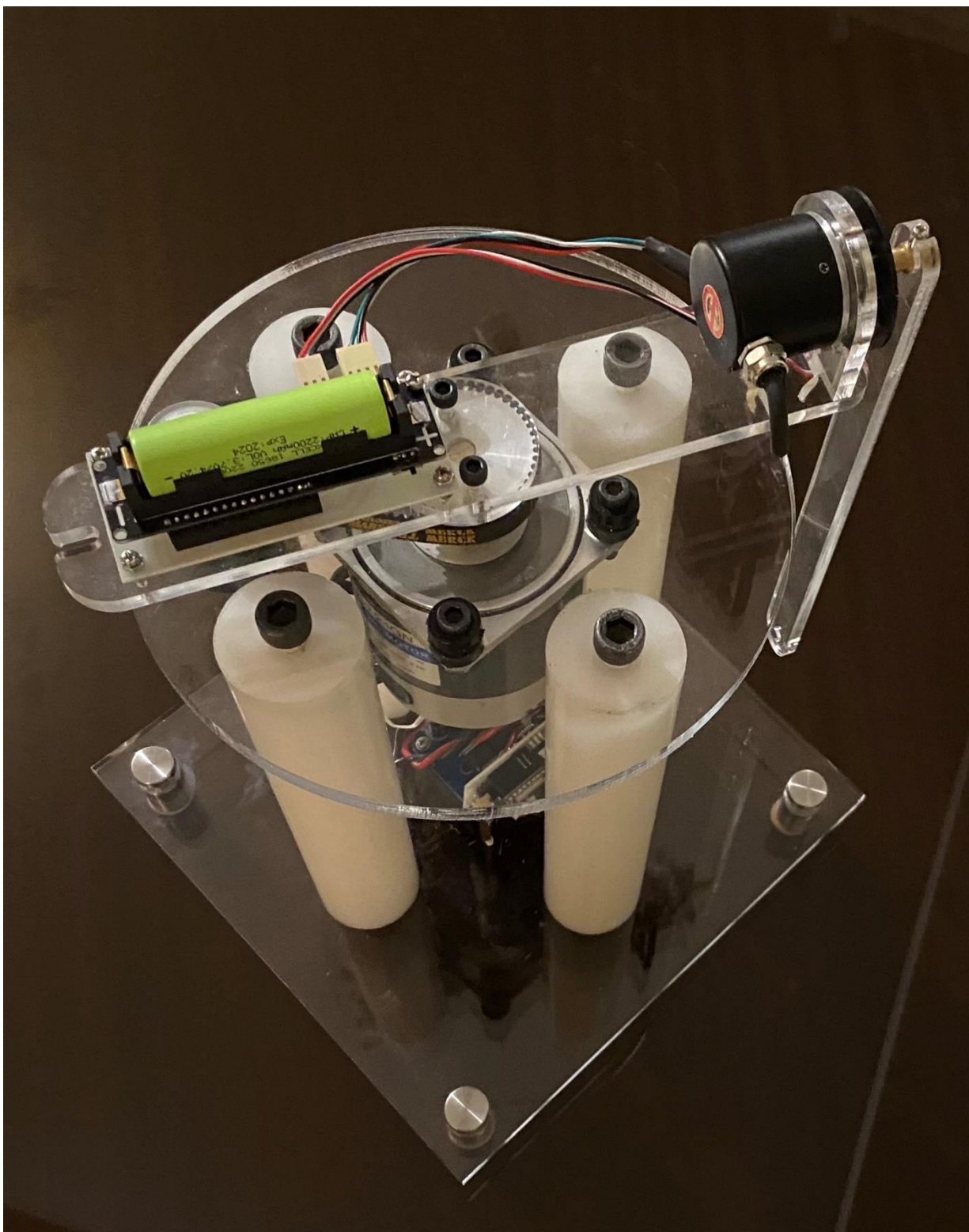
جهت حل مشکل اول، یک محل در انتهای آونگ جهت اضافه کردن وزنه برای بالانس کردن آونگ تعبیه شد. اگرچه به علت چرخش دائمی آونگ رسیدن به بالانس کامل برای سازه ی آونگ و بازوی در حال دوران ممکن نیست، از آنجاکه جرم آونگ نسبت به سایر بخش ها به مراتب کمتر است، امکان رسیدن به بالانس حدودی حول محور دوان وجود دارد. مشکلات دوم و سوم به ترتیب با اضافه کردن پایه ی مناسب به سازه و واشرهای پلاستیکی ضخیم به عنوان دمپر در محل اتصال موتور به بدنه تا حدودی رفع شدند. اما متأسفانه ایراد چهارم همچنان در سازه پابرجا است. همچنین، به علت کهنه گی موتور و بریده شدن بخشی از شفت آن، شفت اندکی از حالت قائم خارج شده است. بنابراین، برای رسیدن به بهترین بازده دستگاه به نظر می رسد بهترین گزینه تعویض موتور با یک موتور جدید است.

با توجه به اعمال تغییرات مذکور در سازه، لازم بود تا آونگ، بازو و نگه دارنده ی انکودر متصل به بازو دوباره طراحی و ساخته شوند. نقشه ی کارگاهی این بخش ها در ضمیمه ی اول آورده شده است.

دو تصویر از سازه ی کامل آونگ معکوس دورانی اصلاح شده در ادامه آورده شده است.



شکل ۲۶: سازه‌ی آونگ معکوس دورانی - نمای روبرو



شکل ۲۷: سازه‌ی آونگ معکوس دورانی - نمای روبرو

۴- محیط شبیه‌سازی^۱

در حین اجرای اصلاحات ذکرشده در بخش ۳ بر روی سازه، سعی شد تا هم‌زمان عملکرد الگوریتم‌های یادگیری تقویتی پیوسته در محیط شبیه‌سازی‌شده در کامپیوتر سنجیده شود. هدف از این کار را به‌طور خلاصه می‌توان در سه مورد زیر بیان کرد:

- ۱- آشنا شدن با نحوه‌ی پیاده‌سازی الگوریتم‌ها
- ۲- یافتن بازه‌های معقول برای پارامترهای اساسی^۲ در الگوریتم‌ها (مانند معماری شبکه‌ی مورد‌استفاده و پارامترهای مورد‌استفاده در آموزش آن از جمله ضریب یادگیری^۳)
- ۳- استفاده از مدل‌های آموزش داده‌شده به کمک شبیه‌ساز^۴ برای کنترل سازه‌ی اصلی

به‌منظور شبیه‌سازی سیستم، در ابتدا یک شبیه‌ساز در زبان برنامه‌نویسی پایتون^۵ از پایه نوشته شد. برای این کار، معادلات دیفرانسیل حاکم بر آونگ معکوس دورانی در حالت گسسته بازنویسی شدند و از آن‌ها برای به‌روزرسانی سینماتیک آونگ معکوس در شبیه‌ساز در زمان‌های گسسته استفاده شد. مهم‌ترین مزیت مدل نوشته‌شده در پایتون، سرعت بالای اجرای آن می‌باشد. در مقابل، با توجه به دینامیک نسبتاً پیچیده‌ی آونگ معکوس دورانی، از بسیاری از جزئیات در شبیه‌سازی اولیه صرف‌نظر شد. در گام دوم شبیه‌سازی، از بخش سیمولینک^۶ نرم‌افزار متلب^۷ برای ساخت یک مدل کامپیوتری دقیق‌تر از آونگ معکوس دورانی استفاده شد. این کار، امکان واردکردن ظریف‌ترین جزئیات در فرایند مدل‌سازی را نیز فراهم نمود اما در عوض، سرعت اجرای مدل پیاده‌سازی شده به نسبت مدل نوشته‌شده در پایتون به‌مراتب پایین‌تر است.

در ادامه، به توضیح مدل دقیق‌تر دو مدل پیاده‌سازی شده می‌پردازیم.

۴-۱- مدل‌سازی در زبان پایتون

این مدل‌سازی پیش از ساخته‌شدن سازه‌ی اولیه‌ی آونگ معکوس دورانی انجام شد. برای مدل‌سازی، آونگ معکوس دورانی به‌صورت یک مدل اجرای محدود^۸ متشکل از یک موتور DC، یک بازوی^۹ متصل به موتور و یک آونگ^{۱۰} در انتهای آن در نظر گرفته شد.

¹ Simulation Environment

² Hyperparameters

³ Learning Rate

⁴ Simulator

⁵ Python

⁶ Simulink

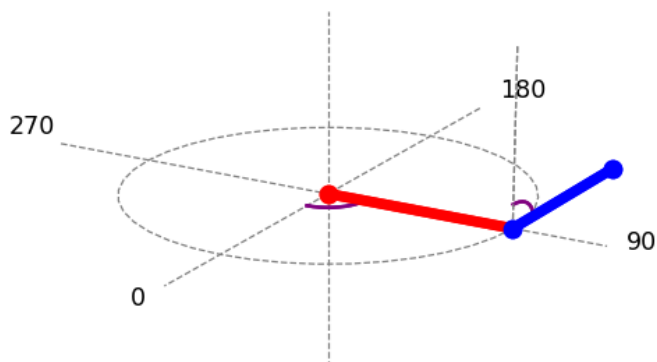
⁷ Matlab

⁸ Lumped

⁹ arm

¹⁰ Pendulum

شکل زیر، نمایی از پنجره گرافیکی نوشته شده به زبان پایتون جهت نمایش وضعیت لحظه‌ای آونگ معکوس دورانی است:



شکل ۲۸: پنجره گرافیکی نوشته شده به زبان پایتون جهت نمایش وضعیت لحظه‌ای آونگ معکوس دورانی شبیه‌سازی شده

در ادامه، به بررسی جنبه‌های مختلف مدل می‌پردازیم.

۴-۱-۱- مفروضات مدل سازی

فرض‌های اساسی ساده کننده اعمال شده در این مدل، به شرح زیر می‌باشند:

- ۱- از وجود تمام اجزای سازه‌ی آونگ معکوس دورانی به جز آونگ و بازو صرف نظر می‌شود.
- ۲- هر دوی بازو و آونگ به صورت دو میله با مقطع دایره‌ای در نظر گرفته می‌شوند.
- ۳- توزیع جرم بازو و آونگ در راستای طولشان یکنواخت است.
- ۴- اصطکاک در محل اتصال آونگ به بازو از نوع اصطکاک ویسکوز^۱ متناسب با سرعت چرخش مفصل در نظر گرفته می‌شود و از سایر انواع اصطکاک صرف نظر می‌شود.
- ۵- از تغییرات مشخصه‌های موتور با افزایش دما صرف نظر می‌شود.

۴-۱-۲- ورودی و خروجی مدل

ورودی و خروجی‌های مدل به شرح زیراند:

جدول ۳: ورودی‌ها و خروجی‌های مدل اولیه

ورودی مدل	
ولتاژ اعمالی به پایه‌های موتور	
خروجی‌های مدل	
مکان زاویه‌ای بازو	مکان زاویه‌ای آونگ
سرعت زاویه‌ای بازو	سرعت زاویه‌ای آونگ
شتاب زاویه‌ای بازو	شتاب زاویه‌ای آونگ

¹ Viscous Friction

۳-۱-۴- پارامترهای مدل سازی

به این ترتیب، پارامترهای دخیل در مدل سازی به شرح موجود در جدول زیر می باشند:

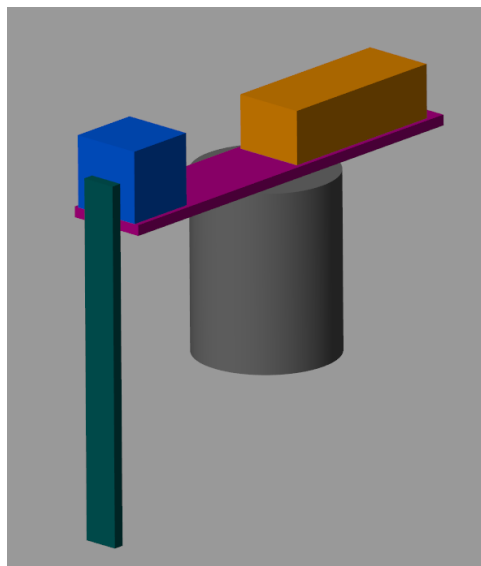
جدول ۴: پارامترهای دخیل در مدل سازی اولیه

پارامترهای دخیل در مدل سازی اولیه	
پارامترهای بازو	چگالی / طول / شعاع مقطع
پارامترهای آونگ	چگالی / طول / شعاع مقطع
مفصل بین آونگ و بازو	ضریب اصطکاک ویسکوز
پارامترهای موتور	مقاومت موتور / اندوکنانس موتور / ضریب Back emf موتور / اینرسی روتور / ضریب میرایی روتور / ضریب گشتاور موتور

در هنگام مدل سازی، به علت در دست نبودن سازه فیزیکی، پارامترهای فوق با مقادیر منطقی جایگزین شدند.

۴-۲- مدل سازی در سیمولینک^۱

در هنگام نوشتن کد مدل اولیه در پایتون، ساخت سازه آغاز نشده بود و جزئیات طراحی نیز به طور کامل مشخص نبود. پس از مشخص شدن جزئیات طراحی سازه، لزوم اعمال این جزئیات در مدل سازی برای رسیدن به اهداف مدل سازی احساس می شد. اما پیاده سازی تمامی جزئیات از پایه به کمک پایتون امری بسیار دشوار و طاقت فرسا بود و به همین دلیل، یک مدل دیگر به کمک بخش سیم اسکپ^۲ در سیمولینک پیاده سازی شد. اگرچه که سرعت اجرای مدل ثانویه به مراتب از مدل اولیه پایین تر است، در آن جزئیات سازه نیز مدل شده اند.



شکل ۲۹: پنجره گرافیکی سمولینک جهت نمایش وضعیت لحظه ای آونگ معکوس دورانی شبیه سازی شده

^۱ Simulink

^۲ Simscape

در ادامه، به بررسی جنبه‌های مختلف این مدل می‌پردازیم.

۴-۲-۱- مفروضات مدل سازی

چند فرض اعمال شده در مدل ثانویه، به شرح زیر می‌باشند:

- ۱- از جرم پیچ‌های موجود در سازه به علت کوچک بودن صرف نظر می‌شود.
- ۲- از جرم مکانیزم کالیبراسیون انکودر متصل به بازو به علت کوچک بودن صرف نظر می‌شود.
- ۲- انکودر و برد چاپی متحرک نصب شده روی بازوی سازه به عنوان اجرام متمرکز در نظر گرفته می‌شود.
- ۳- اصطکاک در محل اتصال آونگ به بازو از نوع اصطکاک ویسکوز^۱ متناسب با سرعت چرخش مفصل در نظر گرفته می‌شود و از سایر انواع اصطکاک صرف نظر می‌شود.
- ۴- از تغییرات مشخصه‌های موتور با افزایش دما صرف نظر می‌شود.

۴-۲-۲- ورودی و خروجی‌های مدل

ورودی مدل نوشته شده در سیمولینک دوره‌ی کاری^۲ سیگنال PWM^۳ ورودی به درایور موتور و خروجی آن، یک سری زمانی از خروجی زمان حال سیستم به همراه تعداد دلخواهی از خروجی‌های آن در زمان‌های گذشته است. به وسیله‌ی این مقادیر و با در دست داشتن زمان نمونه‌برداری^۴ از مدل می‌توان سرعت و شتاب زاویه‌ای آونگ و بازو را محاسبه نمود.

۴-۲-۳- پارامترهای مدل سازی

به این ترتیب، پارامترهای دخیل در مدل سازی به شرح موجود در جدول زیر می‌باشند:

جدول ۵: پارامترهای دخیل در مدل سازی ثانویه

پارامترهای دخیل در مدل سازی دقیق	
پارامترهای بازو	چگالی / طول / عرض و ضخامت مقطع مستطیلی
پارامترهای آونگ	چگالی / طول / عرض و ضخامت مقطع مستطیلی
پارامترهای انکودر	جرم / فاصله مرکز جرم از مرکز دوران / ضریب اصطکاک ویسکوز شفت
پارامترهای مدار چاپی	جرم / فاصله مرکز جرم از مرکز دوران
پارامترهای موتور	مقاومت موتور / اندوکنانس موتور / ضریب Back emf موتور / اینرسی روتور / ضریب میرایی روتور / ضریب گشتاور موتور

¹ Viscous Friction

² Duty Cycle

³ Pulse-Width Modulation

⁴ Sample Time

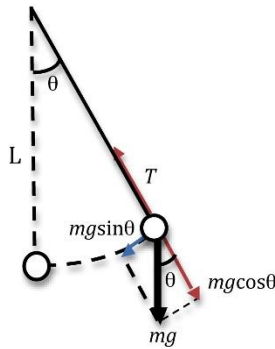
از بین پارامترهای ذکر شده:

- طول، عرض و ضخامت بازو و آونگ جزو پارامترهای طراحی سازه و معلوم هستند.
- جنس بازو و آونگ از پلکسی گلس است و چگالی بازو و آونگ مشخص می باشد.
- جرم انکودر و برد چاپی به کمک ترازو دقیق اندازه گیری شد.
- مراکز جرم انکودر و مدار چاپی برابر مراکز حجم آنها تخمین زده شد و به این ترتیب، فاصله مراکز جرم انکودر و مدار چاپی از مرکز دوران، به دست آمد.

حال، به بررسی روش های به دست آوردن سایر پارامترها می پردازیم.

۴-۲-۱- به دست آوردن ضریب اصطکاک ویسکوز شفت

برای به دست آوردن این ضریب، در ابتدا معادله حاکم بر نوسان آونگ را می نویسم. با نوشتن معادله گشتاور حول نقطه ی دوران برای آونگ نشان داده شده در شکل روبرو خواهیم داشت:



$$I\ddot{\theta} + c\dot{\theta} = mgl\sin\theta \quad (\text{Equation 23})$$

در معادله ی فوق، c ضریب اصطکاک آونگ با مفصل و I اینرسی آونگ حول مرکز دوران است.

شکل ۳۰: آونگ ساده

همان گونه که مشاهده می شود، معادله ی به دست آمده برای نوسان آونگ یک معادله دیفرانسیلی غیرخطی است و حل آن از طریق روابط مرسوم ممکن نیست. اما، در نزدیکی نقطه تعادل آونگ ($\theta \approx 0$) طبق سری مک لورن^۱ برای تابع سینوس، هم عرضی $\sin\theta = \theta$ برقرار است و رابطه ی به دست آمده را می توان به صورت زیر خطی سازی کرد:

$$I\ddot{\theta} + c\dot{\theta} = mgl\theta \quad (\text{Equation 24})$$

معادله ی جدید، یک معادله دیفرانسیل مرتبه دو است و پاسخ آن به شرح زیر است:

$$\theta = c_1 e^{-\frac{\sqrt{c^2 + 4Imgl} + c}{2I}x} + c_2 e^{\frac{\sqrt{c^2 + 4Imgl} - c}{2I}x} \quad (\text{Equation 25})$$

¹ Maclaurin Series

در حل فوق، مقادیر c_1 و c_2 از شرایط مرزی برای مکان و سرعت زاویه‌ای آونگ به دست می‌آیند. اگر آونگ از زاویه اولیه‌ی θ_0 با مقدار نزدیک به صفر و بدون سرعت اولیه رها شود، می‌توان ضرایب c_1 و c_2 را به صورت زیر محاسبه نمود:

$$c_1 = \frac{\sqrt{c^2 + 4Imgl} - c}{2\sqrt{c^2 + 4Imgl}} \quad , \quad c_2 = \frac{\sqrt{c^2 + 4Imgl} + c}{2\sqrt{c^2 + 4Imgl}}$$

با جاگذاری در معادله اصلی خواهیم داشت:

$$\theta = \frac{\sqrt{c^2 + 4Imgl} - c}{2\sqrt{c^2 + 4Imgl}} \theta_0 e^{-\frac{\sqrt{c^2 + 4Imgl} + c}{2I}x} + \frac{\sqrt{c^2 + 4Imgl} + c}{2\sqrt{c^2 + 4Imgl}} \theta_0 e^{\frac{\sqrt{c^2 + 4Imgl} - c}{2I}x} \quad (\text{Equation 25})$$

برای اندازه‌گیری ضریب اصطکاک ویسکوز در شفت انکودر (اتصال بین بازو و آونگ) از معادله‌ی فوق استفاده شد. برای آونگ معکوس دورانی ساخته شده، تمامی پارامترهای دخیل در معادله فوق، به جز مقدار C مشخص است. پس با رها کردن آونگ از یک زاویه‌ی θ_0 نزدیک به صفر، و مشاهده‌ی رفتار آونگ، می‌توان مقدار ضریب اصطکاک ویسکوز را محاسبه کرد.

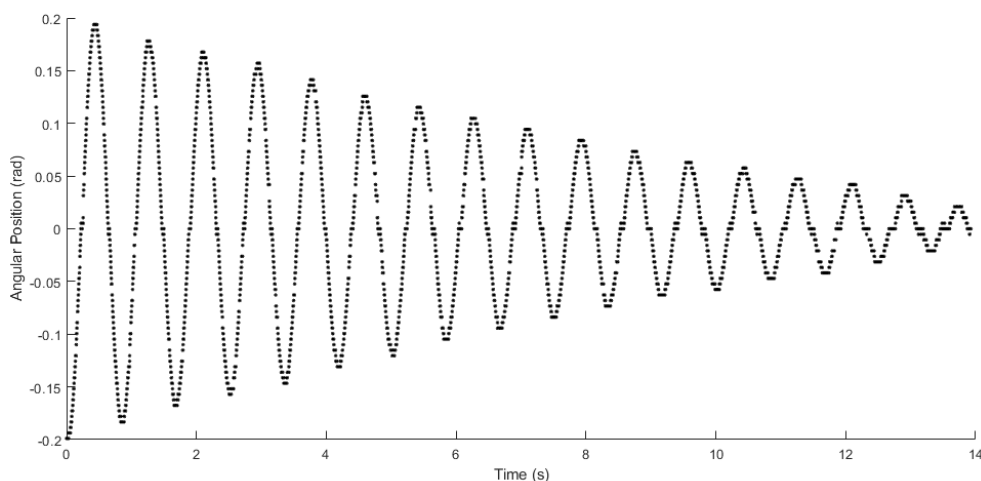
به این ترتیب، مراحل محاسبه‌ی ضریب اصطکاک ویسکوز به صورت زیر خواهد بود:

گام ۱: آونگ را از نقطه تعادل به اندازه‌ی کم منحرف می‌کنیم به گونه‌ای که فرض خطی ساز معادله دیفرانسیل آونگ ($\sin\theta = \theta$) تقریباً برقرار شود.

گام ۲: نوسانات آونگ تا زمان سکون را به کمک انکودر اندازه‌گیری می‌کنیم.

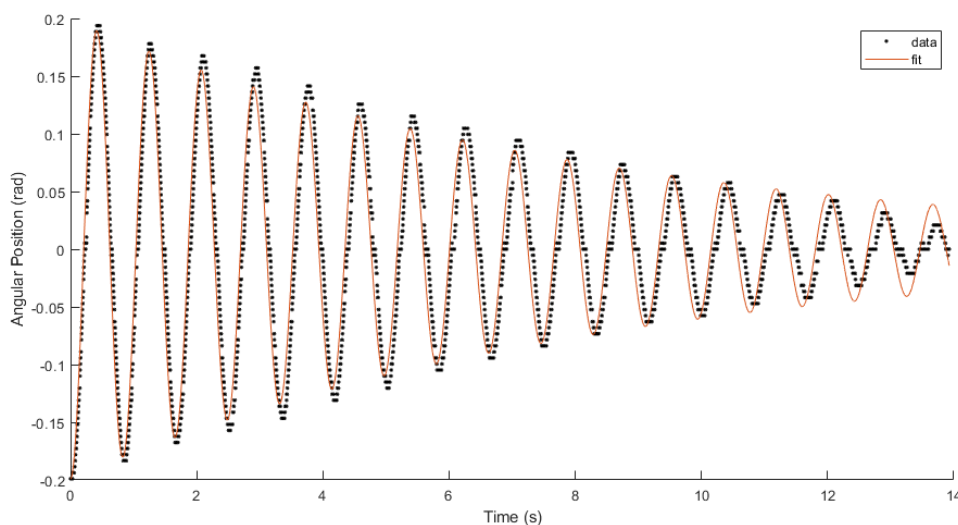
گام ۳: به کمک برازش منحنی به فرم جواب به دست آمده، مقدار پارامتر C را تعیین می‌نماییم.

منحنی مکان زاویه‌ای آونگ از زمان رها شدن تا هنگام سکون در زیر آورده شده است:



شکل ۳۱: نمودار مکان زاویه‌ای بر حسب زمان برای آونگ در هنگام رها شدن از زاویه θ_0

حال، با در دست داشتن پارامترهای آونگ و به کمک Tool Box برازش منحنی نرم‌افزار متلب یک منحنی به فرم جواب به‌دست‌آمده و با پارامتر C مجهول به داده‌های به‌دست‌آمده از آزمایش برازش می‌کنیم. منحنی برازش شده در ادامه آمده است.



شکل ۳۲: برازش منحنی نمودار مکان زاویه‌ای بر حسب زمان برای آونگ در هنگام رها شدن از زاویه θ_0

همچنین، مقدار به دست آمده برای C از طریق برازش منحنی در زیر آورده شده است:

$$c = 0.000211 \text{ N.m.s/rad } (0.0002093, 0.0002127) \text{ 95\% confidence}$$

۴-۲-۲- به دست آوردن پارامترهای موتور (مکانیزم اندازه گیری مکان آن)

لازم است تا تابع تبدیل موتور DC را به دست آورده‌یم. برای این کار، ابتدا معادلات موتور از بخش ۲-۲ را بار دیگر بازنویسی می‌کنیم:

$$R_a i_a + L_a \frac{di_a}{dt} + V_b = V_a \quad (\text{Equation 19}) \quad V_b = K_b \frac{d\theta}{dt} \quad (\text{Equation 20})$$

$$I \frac{d^2\theta}{dt^2} + C \frac{d\theta}{dt} = T - T_L \quad V_b = K_b \frac{d\theta}{dt} \quad (\text{Equation 21}) \quad T = K_t i_a \quad (\text{Equation 22})$$

متأسفانه، باوجود جست‌وجوی فراوان، کاتالوگ موتور DC مورد استفاده یافت نشد. لذا لازم بود تا این پارامترهای دخیل در موتور از روی رفتار فیزیکی آن اندازه‌گیری شوند.

در طی این فرایند، انکودر به وسیله‌ی تسمه و پولی همواره به موتور DC متصل بود و این دو بخش سیستم، به عنوان یک بخش واحد در نظر گرفته شدند. این موضوع، موجب می‌شود که اینرسی پولی‌ها و تسمه و اصطکاک ویسکوز در انکودر نیز در محاسبات به طور خودکار در محاسبات وارد شوند.

در ابتدا، به کمک یک RLC متر، مقادیر مقاومت الکتریکی^۱ بین پایه‌های موتور و اندوکدانس^۲ اندازه‌گیری شد. مقادیر اندازه‌گیری شده به شرح زیر می‌باشند:

$$R_a = 2.5 \, \Omega \quad , \quad L_a = 2.1 \, mH$$

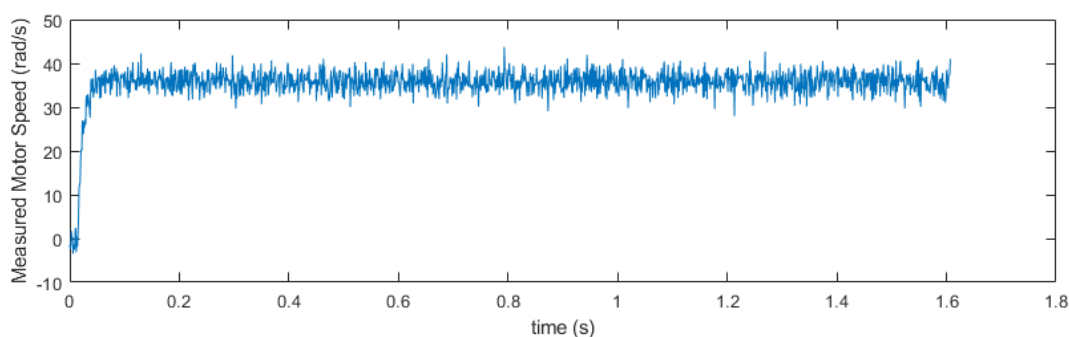
لازم به ذکر است که در فرایند مدل سازی از هرگونه تغییرات این مقادیر برحسب دما صرف نظر شد.

در ادامه‌ی کار، جهت به دست آوردن ضریب $Back \, emf$ (K_b) مراحل زیر را دنبال می‌کنیم:

- گام ۱: شفت موتور را به سر دریل متصل می‌کنیم و آن را می‌چرخانیم.
- گام ۲: ولتاژ پایه‌های موتور (همان ولتاژ $Back \, emf$) را به کمک پین آنالوگ آردوینو^۳ و سرعت دوران شفت موتور را به کمک انکودر برای مدتی اندازه‌گیری می‌کنیم.
- گام ۳: از معادله‌ی ۲، می‌توان K_b را به صورت زیر حساب کرد:

$$V_b = K_b \frac{d\theta}{dt} \rightarrow K_b = \frac{V_b}{\dot{\theta}}$$

نمودار اندازه‌ی سرعت دورانی موتور، بر حسب ولتاژ اندازه‌گیری شده در پایه‌های آن در ادامه آورده شده است:

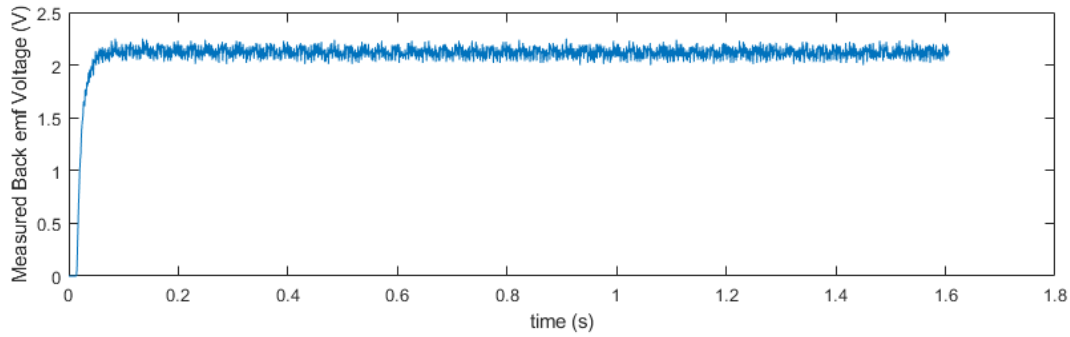


شکل ۳۳: نمودار سرعت دوران شفت موتور در هنگام چرخاندن آن به وسیله دریل

¹ Resistance

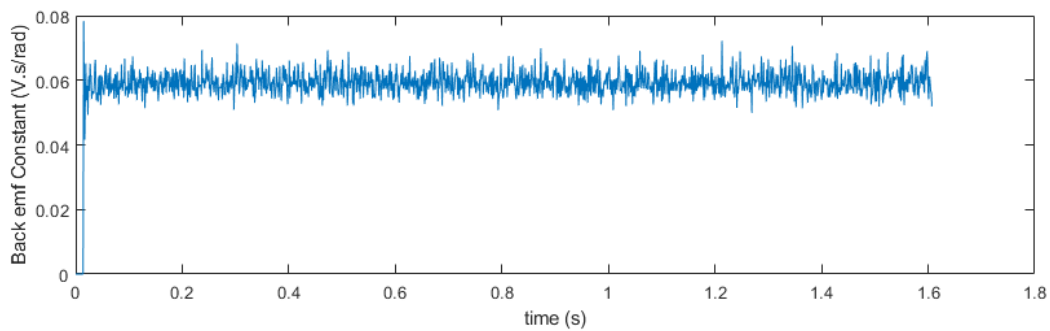
² Inductance

³ Arduino



شکل ۳۴: نمودار ولتاژ اندازه‌گیری شده در پایه‌های موتور در هنگام چرخاندن شفت آن به وسیله دریل

همچنین، مقادیر متناظر K_b حاصل از تقسیم دو نمودار قبلی بر هم در نمودار زیر آورده شده است:



شکل ۳۵ مقادیر به‌دست آمده برای K_b در طول زمان

به این ترتیب، مقدار متوسط K_b به صورت زیر به دست می‌آید:

$$K_b = 0.0587 \frac{Vs}{rad}$$

تا به اینجا، مقادیر R_a ، L_a و k_b محاسبه شدند. برای به دست آوردن سایر پارامترهای موتور، از پاسخ پله‌ی^۱ آن استفاده شد. برای این کار، لازم است تا تابع تبدیل^۲ موتور محاسبه شود. برای محاسبه‌ی تابع تبدیل، در ابتدا از طرفین معادلات موتور تبدیل لاپلاس^۳ می‌گیریم:

$$R_a i_a(s) + L_a s i_a(s) + V_b(s) = V_a(s) \quad (\text{Equation 26}) \quad V_b(s) = K_b s \theta(s) \quad (\text{Equation 27})$$

$$I s^2 \theta(s) + C s \theta(s) = T(s) - T_L \quad (\text{Equation 28}) \quad T(s) = K_t i_a(s) \quad (\text{Equation 29})$$

¹ Step Response

² Transfer Function

^۳ Laplace Transform

با جاگذاری مقدار $V_b(s)$ از معادله ۲۷ در معادله ۲۶ و $T(s)$ از معادله ۲۹ در معادله ۲۸ به دست می آید:

$$R_a i_a(s) + L_a s i_a(s) + K_b s \theta(s) = V_a(s) \quad (\text{Equation 30})$$

$$I s^2 \theta(s) + C s \theta(s) = K_t i_a(s) - T_L \quad (\text{Equation 31})$$

همچنین، با جاگذاری مقدار $i_a(s)$ از معادله ۳۰ در معادله ۳۱ و اعمال $T_L = 0$ ، خواهیم داشت:

$$\frac{\theta(s)}{V_a(s)} = \frac{k_t}{(R_a + L_a s)(I s^2 + c s) + k_b k_t s}$$

نهایتاً داریم:

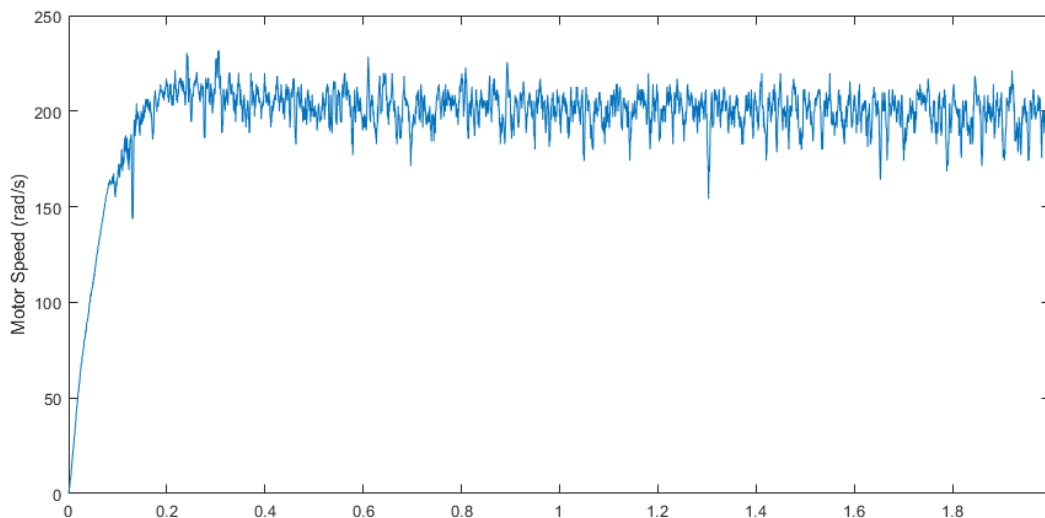
$$G_1 = \frac{\theta(s)}{V_a(s)} = \frac{k_t}{I L_a s^3 + (I R_a + c L_a) s^2 + (R_a c + k_b k_t) s} \quad (\text{Equation 32})$$

همچنین، اگر خروجی تابع تبدیل را سرعت زاویه‌ای در نظر بگیریم:

$$G_1 = \frac{\dot{\theta}(s)}{V_a(s)} = \frac{k_t}{I L_a s^2 + (I R_a + c L_a) s + (R_a c + k_b k_t)} \quad (\text{Equation 32})$$

در تابع تبدیل به دست آمده مقادیر I ، c و k_t مجهول هستند. با به دست آوردن ضرایب تابع تبدیلی به فرم به دست آمده از پاسخ پله^۱ موتور می توان پارامترهای مجهول را پیدا کرد.

پاسخ پله‌ی موتور به ازای ولتاژ 24V و با فرکانس نمونه برداری^۲ ۲ کیلوهرتز به مانند زیر است:

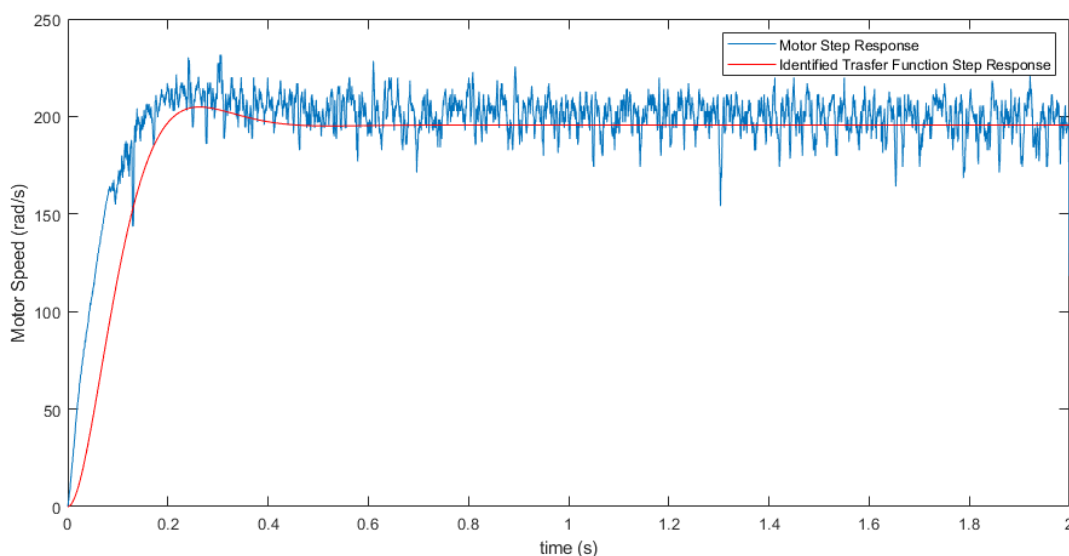


شکل ۳۶: پاسخ پله‌ی موتور به ازای ولتاژ 24V

¹ Step Response

² Sampling Frequency

همچنین، پاسخ پله‌ی تابع تبدیل به‌دست‌آمده برای پاسخ فوق، در toolbox شناسایی سیستم^۱ متلب به شرح زیر است:



شکل ۳۷: پاسخ پله‌ی موتور در کنار پاسخ پله‌ی تابع تبدیل شناسایی‌شده برای آن

فرم تابع تبدیل شناسایی‌شده به‌صورت زیر است:

$$\hat{G}_1 = \frac{\dot{\theta}(s)}{V_a(s)} = \frac{1}{0.0004432s^2 + 0.01027s + 0.1227}$$

با مقایسه‌ی این تابع با G_1 ، مجهولات موجود در معادلات به‌دست می‌آید:

$$\frac{IL_a}{k_t} = 0.0004432, \quad \frac{IR_a + cL_a}{k_t} = 0.01027, \quad \frac{R_ac + k_bk_t}{k_t} = 0.1227$$

با معلوم بودن R_a ، L_a و k_b ، دستگاه سه معادله و سه مجهول غیرخطی فوق به‌راحتی حل می‌شود و مقادیر I ، K_t و c به‌صورت زیر به دست می‌آیند:

$$k_t = 0.0584 \text{ N.m/A}, \quad I = 0.0002 \text{ kg.m}^2, \quad c = 0.0015 \text{ N.m.s/rad}$$

به‌این ترتیب، فرایند شناسایی پارامترهای موتور به پایان می‌رسد.

¹ System Identification

۴-۲-۳- به دست آوردن مشخصه انتقالی درایور موتور^۱

برای به حرکت درآوردن موتور، از یک درایور موتور مدل BTS7960 استفاده می‌شود. ورودی درایور، یک سیگنال PWM^۲ است و متناسب با دوره کاری^۳ این سیگنال، ولتاژ اعمالی در پایانه‌های موتور تعیین می‌شود. ولتاژ خروجی درایور موتور را می‌توان با تقریب مناسب به صورت خطی نسبت به دوره کاری سیگنال PWM ورودی به آن در نظر گرفت. در تعیین مشخصه انتقالی درایور، موارد زیر حائز اهمیت است:

- ۱- افت ولتاژ درونی درایور مطابق مقدار عنوان شده در کاتالوگ و مقادیر مشاهده شده در عمل برای سازه در حدود $0.5V$ می‌باشد. این مقدار تقریباً مستقل از شدت جریان گرفته شده از درایور است.
- ۲- با توجه به این که بیشینه ولتاژ خروجی درایور برابر ولتاژ منبع تغذیه ($25V$) و در حالت دوره کاری 100% است، شیب تغییرات ولتاژ خروجی درایور به دوره کاری با تقریب مناسب برابر 0.25 است.
- ۳- با تست کردن روی سازه‌ی آونگ معکوس، آستانه‌ی حرکت پیوسته‌ی دستگاه حدوداً در ولتاژ $4.5V$ بین دو پایه موتور است.

با توجه به موارد فوق، می‌توان مشخصه انتقالی درایور موتور را به صورت زیر به دست آورد:

$$V_a = \text{MAX}(0.25 \times \text{Duty Cycle} - 4.5, 0)$$

در عبارت فوق، V_a ، ولتاژ خروجی درایور و پایه‌های موتور و Duty Cycle بر حسب درصد بیان می‌شود.

به این ترتیب، مقدار تمامی پارامترهای دخیل در مدل به دست می‌آیند. بار دیگر، مقدار تمامی این پارامترها در جداول زیر آورده شده است:

جدول ۶: مقادیر پارامترهای دخیل در مدل سازی بازو

پارامترهای دخیل در مدل سازی بازو	
پارامتر	مقدار
چگالی	1180 kg/m^3
طول	270 mm
عرض	47 mm
ضخامت	6 mm

¹ Motor Driver

² Pulse-Width Modulation

³ Duty Cycle

جدول ۷: مقادیر پارامترهای دخیل در مدل‌سازی آونگ

پارامترهای دخیل در مدل‌سازی آونگ	
پارامتر	مقدار
چگالی	1180 kg/m^3
طول	250 mm
عرض	20 mm
ضخامت	6 mm

جدول ۸: مقادیر پارامترهای دخیل در مدل‌سازی موتور DC (و انکودر متصل به آن)

پارامترهای دخیل در مدل‌سازی موتور DC (و انکودر متصل به آن)	
پارامتر	مقدار
مقاومت موتور	2.5Ω
اندوکنانس موتور	2.1 mH
ضریب Back emf موتور	0.0587 V.s/rad
اینرسی روتور	0.0002 kg.m^2
ضریب میرایی روتور	0.0015 N.m.s/rad
ضریب گشتاور موتور	0.0584 N.m/A

جدول ۹: مقادیر پارامترهای دخیل در مدل‌سازی انکودر متصل به آونگ

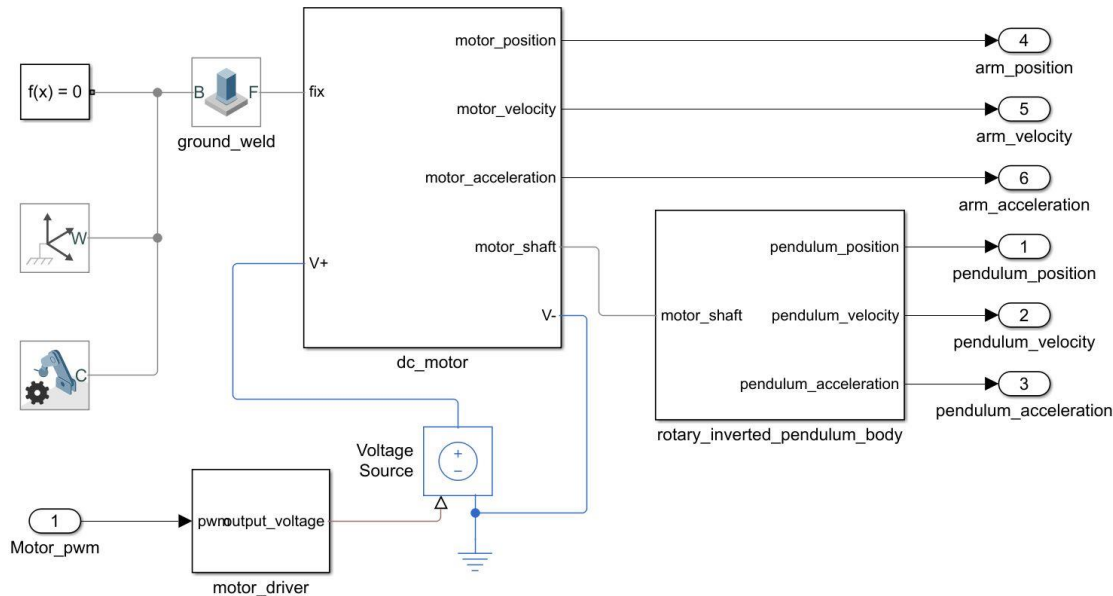
پارامترهای دخیل در مدل‌سازی انکودر متصل به آونگ	
پارامتر	مقدار
جرم	95 g
فاصله مرکز جرم از مرکز دوران	97 mm
ضریب اصطکاک ویسکوز شفت	0.0002 N.m.s/rad
ضخامت	6 mm

جدول ۱۰: مقادیر پارامترهای دخیل در مدل‌سازی بورد چاپی متصل به بازو

پارامترهای دخیل در مدل‌سازی بورد چاپی متصل به بازو	
پارامتر	مقدار
جرم	100 g
فاصله مرکز جرم از مرکز دوران	50 mm

۴-۲-۳- پیاده‌سازی سیستم در سیمولینک^۱

برای پیاده‌سازی سیستم آونگ معکوس دورانی در سیمولینک، از بخش Simscape استفاده شد. این بخش امکان شبیه‌سازی انواع سیستم‌های فیزیکی اعم از مغناطیسی، حرارتی، دینامیکی، ارتعاشاتی، الکتریکی و و از همه مهم‌تر امکان ایجاد ارتباطات بین این سیستم‌ها را برای کاربر مهیا می‌کند. در زیر، شمای کلی پیاده‌سازی مدل پیاده‌سازی شده آورده شده است:



شکل ۳۸: شمای کلی سیستم آونگ معکوس دورانی در سیمولینک

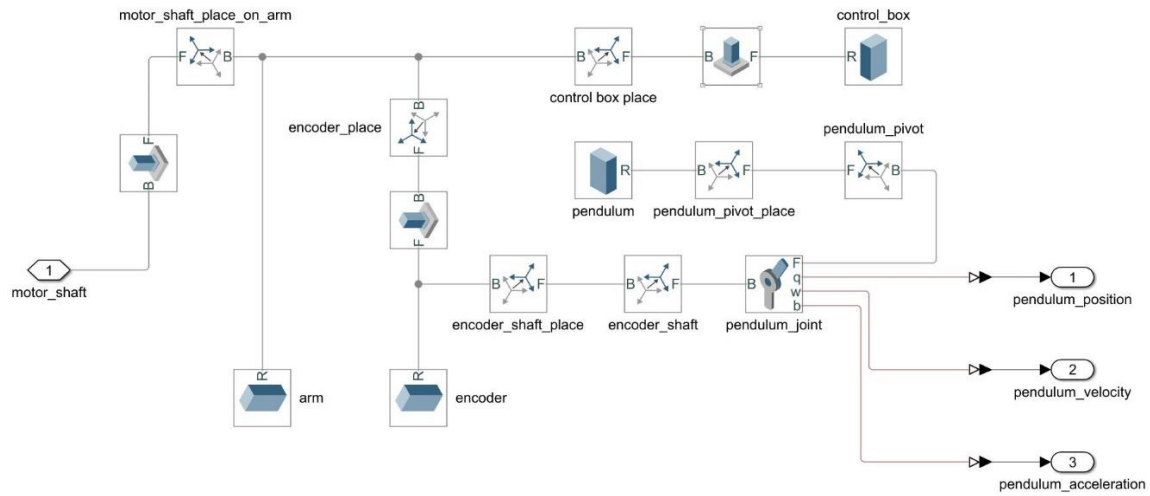
در تصویر فوق، سه زیرسیستم^۲ دیده می‌شود:

- ۱- زیرسیستم `rotary_inverted_pendulum_body` که در آن بخش بازو و آونگ سیستم آونگ معکوس دورانی به کمک بلوک‌های بخش `Multibody` از بخش `Simscape` پیاده‌سازی شده‌اند
- ۲- زیرسیستم `dc_motor` که در آن مدلی از موتور `dc` پیاده‌سازی شده است.
- ۳- زیرسیستم `motor_driver` که یک تابع در زبان برنامه‌نویسی `Matlab` است که ورودی آن مقدار سیگنال `pwm` ورودی به درایور و ولتاژ اعمالی بر پایه‌های موتور است.

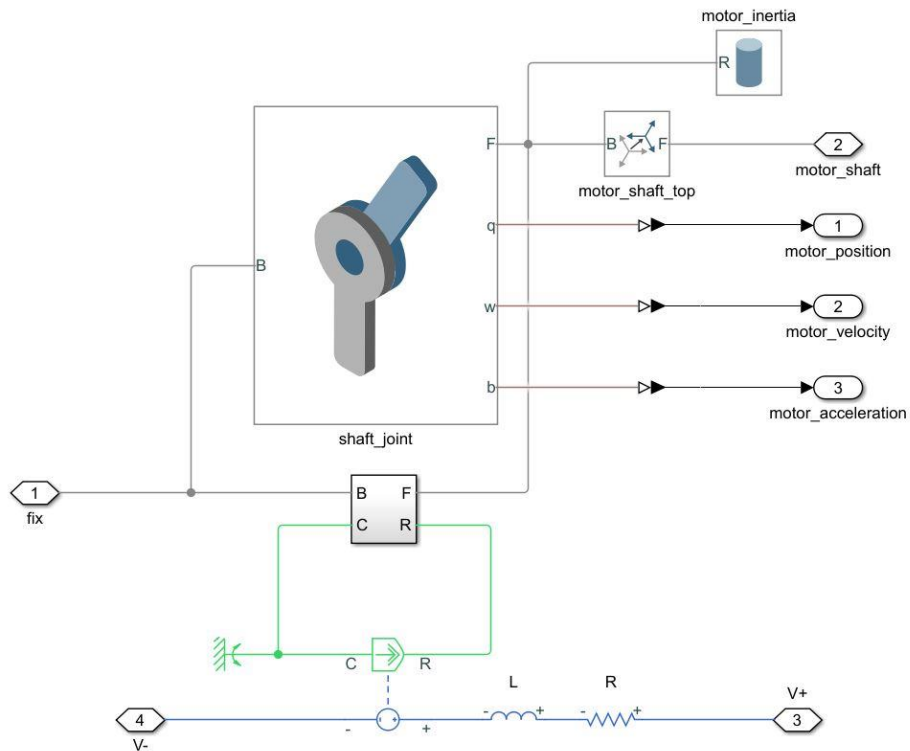
^۱ Simulink

^۲ Subsystem

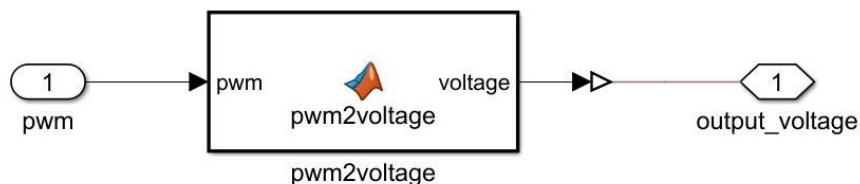
در سه شکل بعدی، نمایی از بخش‌های این سه زیرسیستم آورده شده است:



شکل ۳۹: نمایی از زیرسیستم *rotary_inverted_pendulum_body*



شکل ۴۰: نمایی از زیرسیستم *dc_motor*



شکل ۴۱: نمایی از زیرسیستم *motor_driver*

لازم به ذکر است که پارامترهای اساسی مدل فوق تماماً به صورت متغیر^۱ در محیط متلب پیاده سازی شده اند تا فرایند تغییر و به روز رسانی آن ها سهولت یابد.

۴-۲-۴- کار با محیط شبیه سازی

کار با محیط شبیه سازی گرافیکی در سیمولینک از دو طریق میسر است:

- ۱- به طور مستقل از محیط متلب
- ۲- به طور غیرمستقل از محیط پایتون

نرم افزار متلب بستری مناسب برای پیاده سازی انواع الگوریتم های کنترلی کلاسیک است و برای پیاده سازی این روش ها می توان از امکانات نرم افزار متلب به همراه مدل سیمولینک آونگ معکوس دورانی استفاده کرد.

از طرف دیگر، پایتون، قوی ترین زبان برنامه نویسی جهت پیاده سازی الگوریتم های یادگیری ماشینی^۲ می باشد. به کمک افزونه^۳ *Matlab engine* در نرم افزار متلب و کتابخانه ی *engine.matlab* در زبان پایتون، می توان یک درگاه ارتباطی بین متلب و پایتون ایجاد کرد و به کمک این درگاه می توان مدل آونگ معکوس دورانی را کنترل کرد. جهت برقراری این ارتباط و کنترل مدل، یک کلاس به نام *SimulinkModel* در زبان پایتون نوشته شد که با فراخوانی آن می توان به مدل سیمولینک متصل شد.

¹ Variable

² Machine Learning

³ Plug-in

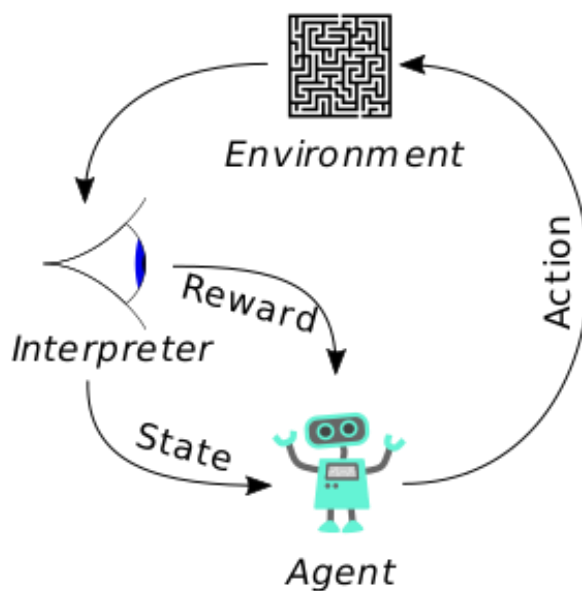
۵- یادگیری تقویتی^۱

یادگیری تقویتی یکی از شاخه‌های یادگیری ماشین^۲ است که با الهام گرفتن از روند طبیعی یادگیری موجودات زنده، به کشف شیوه‌ی عملکرد صحیح در یک محیط^۳ می‌پردازد. در طی فرایند فراگیری رفتار صحیح که اصطلاحاً یادگیری^۴ نامیده می‌شود، یک یا چند عملگر^۵ از طریق تعامل با یک محیط سعی در بیشینه کردن پاداش^۶ دریافتی خود از محیط دارند. یادگیری تقویتی، در کنار یادگیری با نظارت و یادگیری بدون نظرات، از انواع الگوهای یادگیری در حوزه یادگیری ماشین می‌باشد.

در این بخش، به بررسی مفاهیم مرتبط با حوزه یادگیری تقویتی خواهیم پرداخت.

۵-۱- تعاریف پایه

شماتیک زیر، کلیت فرایند طی شده در الگوریتم‌های یادگیری تقویتی را نشان می‌دهد:



شکل ۴۲: شماتیک فرایندهای یادگیری تقویتی^۷

¹ Reinforcement Learning

² Machine Learning

³ Environment

⁴ Learning

⁵ Agent

⁶ Reward

^۷ برگرفته از Wikipedia.org

حال، به بررسی مفاهیم کلیدی در سازوکار نشان داده شده در شکل می پردازیم:

- **عملگر (Agent):** تصمیم گیرنده در محیط است که سعی دارد از طریق بیشینه کردن پاداش دریافتی اش، رفتار صحیح در محیط را یاد بگیرد.
- **محیط (Environment):** سازوکاری است که Agent با آن در تعامل است و سعی در یادگیری عملکرد بهینه در آن دارد.
- **وضعیت (State):** حالات ممکن از محیط است که Agent می تواند در آن ها قرار گیرد.
- **عمل (Action):** کاری است که Agent می تواند متناسب با یک State خاص انجام دهد.
- **پاداش (Reward):** جایزه ای است که Agent متناسب با عملکردش از محیط دریافت می نماید. عملکردهای بهتر، پاداش های بیشتری دریافت می کنند.

سایر مفاهیم اساسی یادگیری تقویتی به شرح زیر هستند:

- **سیاست (Policy):** استراتژی است که Agent آن را در محیط دنبال می کند و با توجه به آن در State های مختلف، Action های متناسب را انتخاب می نماید. به عبارت دیگر، Policy تابعی است که State یک Agent را به یک Action نگاشت می کند.
- **سیاست بهینه (Optimal Policy):** سیاستی است که پاداش ناشی از دنبال کردن آن، از پاداش ناشی از دنبال کردن همه ی سیاست های ممکن دیگر بیشتر باشد. ثابت می شود که در یک مسئله ی یادگیری تقویتی، همواره سیاست بهینه به طور یکتا وجود دارد.
- **اپیزود (Episode):** مجموعه اعمال یک Agent از زمان شروع فرایند تصمیم گیری در محیط تا پایان را یک اپیزود می نامند. فرایندهای یادگیری تقویتی که در قالب چند اپیزود انجام می شوند را اپیزودیک^۱ و سایر فرایندها را ادامه دار^۲ می نامند. در بسیاری از مواقع، با تعریف یک نقطه ی شروع و پایان دلخواه برای فرایندهای ادامه دار، آن ها را به فرایندهای اپیزودیک تبدیل می کنیم.
- **بهره برداری از دانش قبلی (Exploitation):** زمانی است که Agent از دانش به دست آمده ی خود در محیط برای تصمیم گیری استفاده می کند.
- **اکتشاف (Exploration):** زمانی است که Agent بدون توجه به دانش به دست آمده ی خود در محیط، برای تصمیم گیری دست به اعمال رندم می زند.

¹ Episodic Tasks

² Continuing Tasks

- استفاده از دانش قبلی در مقابل اکتشاف (Exploration Vs. Exploitation): هردوی

تصمیم‌گیری بر مبنای دانش پیشین و یا اکتشاف، مزایا و معایب منحصر به خود را دارند. در هنگامی که Agent دانش محیطی کافی جمع‌آوری کرده باشد، تکیه بر دانشش عملی عاقلانه به نظر می‌رسد. اما اگر تمام تصمیم‌گیری‌های Agent بر مبنای دانش پیشینش باشد، Agent ممکن است بسیاری از Action ها را تجربه نکرده باشد و نتواند دانش خود را در مورد رفتار صحیح در محیط و بیشینه کردن پاداش تکمیل نماید. از طرفی، استفاده از دانش قبلی در هنگام فرایند یادگیری عملاً بی‌معنا است. این موضوع را اصطلاحاً Exploration Vs. Exploitation Trade-off می‌نامند. روش‌های زیادی جهت رفع این مشکل معرفی شده است که از میان آن‌ها می‌توان به سیاست epsilon-greedy اشاره نمود. در این سیاست، Agent در هر مرحله تصمیم‌گیری به احتمال ϵ عمل تصادفی انجام می‌دهد و اکتشاف می‌کند و به احتمال $1 - \epsilon$ بر دانش خود تکیه می‌کند. لازم به ذکر است که ϵ یک عدد مثبت بین صفر و یک می‌باشد. مرسوم است که باگذشت زمان و افزایش دانش Agent مقدار ϵ را کاهش می‌دهند تا جایی که نهایتاً Agent برای تصمیم‌گیری تنها به دانش پیشین خود متوسل می‌شود.

نمادگذاری متداول برای مفاهیم فوق در ادبیات حوزه‌ی یادگیری تقویتی در جدول زیر نمایش داده شده‌اند.

جدول ۱۱: نمادگذاری متداول برای مفاهیم پایه‌ای در یادگیری تقویتی

نماد	مفهوم
S	<i>State</i>
A	<i>Action</i>
R	<i>Reward</i>
π	<i>Policy</i>
π_*	<i>Optimal Policy</i>
E	<i>Episode</i>

اندیس t در زیر هر یک از نمادهای S ، A و R به معنی مربوط بودن نماد به زمان t ام اپیزود است. به عنوان مثال، S_t به معنی State بازدید شده در زمان t است.

۵-۲- هدف نهایی مسائل یادگیری تقویتی

در مسائل حوزه‌ی یادگیری تقویتی، هدف نهایی، فراگیری یک سیاست (Policy) بهینه برای تصمیم‌گیری در محیط است که Agent به کمک آن بتواند Reward دریافتی خود از محیط را بیشینه سازد. به عبارت دیگر، در طی فرایند آموزش، Agent با استفاده از مشاهدات محیطی سعی می‌کند سیاست خود را به سیاست بهینه (π_*) نزدیک نماید.

۵-۳- مسائل یادگیری تقویتی گسسته

در مسائل یادگیری تقویتی گسسته، State ها و Action ها دارای مقادیر گسسته هستند. در مقابل مسائل یادگیری تقویتی گسسته، آن دسته از مسائل که State ها و Action هایشان مقادیر پیوسته دارد.

در ادامه، به بررسی فرمولاسیون مسائل یادگیری تقویتی پیوسته و شیوهی حل آن‌ها خواهیم پرداخت.

۵-۳-۱- فرمولاسیون ریاضی مسائل یادگیری گسسته

همان‌طور که پیش‌تر گفته شد، در مسائل یادگیری تقویتی گسسته، مقادیر State و Action همواره گسسته است. برای مدل‌سازی این دسته از مسائل، از فرایند تصمیم‌گیری مارکوف^۱ استفاده می‌شود.

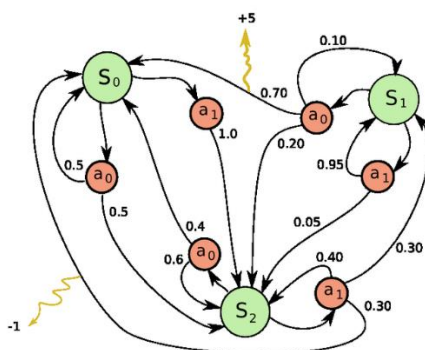
۵-۳-۱-۱- فرایند تصمیم‌گیری مارکوف

فرایند تصمیم‌گیری مارکوف چارچوبی ریاضی است برای مدل‌سازی تصمیم‌گیری در شرایطی که نتایج تا حدودی تصادفی و تا حدودی تحت کنترل یک تصمیم‌گیر است. فرایندهای تصمیم‌گیری مارکوف، فرایندهای کنترلی تصادفی زمان گسسته هستند. در این فرایندها، تصمیم‌گیرنده در زمان t با اتخاذ تصمیم A_t به احتمال P از وضعیت S_t به وضعیت S_{t+1} منتقل می‌شود و پاداش R_{t+1} را از محیط دریافت می‌کند.

$$S_t \rightarrow A_t \rightarrow R_{t+1} \rightarrow S_{t+1}$$

این روند، تا زمان رسیدن به State پایانی^۲ ادامه می‌یابد.

شکل زیر، یک فرایند تصمیم‌گیری مارکوف با سه حالت (State) و دو عمل (Action) را نشان می‌دهد. احتمال‌های انتقال از هر وضعیت به وضعیت دیگر در شکل نشان داده شده‌اند.



شکل ۴۳: یک فرایند تصمیم‌گیری مارکوف ساده با سه حالت (State) و دو عمل (Action)^۳

^۱ Markov Decision Process

^۲ Terminal State

^۳ برگرفته از Wikipedia.org

فرایندهای تصمیم‌گیری مارکوف، به علت نزدیکی به مفاهیم پایه‌ای مسائل یادگیری تقویتی گسسته، بستری مناسب برای مدل‌سازی این مسائل می‌باشند.

۵-۳-۱-۲- بهینه‌سازی در تصمیم‌گیری مارکوف

در فرایندهای تصمیم‌گیری، قصد داریم سلسله تصمیماتی^۱ را اتخاذ نماییم که منجر به بیشترین پاداش^۲ دریافتی از محیط شود. به‌طور مشخص، در زمانی مشخص t ، قصد داریم به‌گونه‌ای عمل کنیم که عبارت زیر بیشینه شود:

$$G_t = E \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \right]$$

در عبارت فوق، E امید ریاضی است و به علت وجود عدم قطعیت در فرایند تصمیم‌گیری مارکوف وارد معادلات می‌شود. γ نیز نرخ کاهش^۳ نام دارد و یک عدد بین صفر و یک است. مقدار γ ، میزان توجه به پاداش آینده را مشخص می‌کند. به ازای $\gamma = 0$ ، Agent تمام تصمیم‌گیری‌ها را با توجه به پاداش سریع و به‌صورت حریصانه^۴ اتخاذ می‌کند. در مقابل، به ازای $\gamma = 1$ ، Agent در هنگام تصمیم‌گیری نه‌تنها مقادیر پاداش فوری، بلکه تمامی پاداش‌های آینده را نیز در نظر می‌گیرد. این ضریب یکی از پارامترهای اساسی^۵ در الگوریتم‌های یادگیری تقویتی است و بسته به شرایط مقدار آن انتخاب می‌شود.

۵-۳-۱-۳- مفهوم State Value و Action Value

در راستای بیشینه کردن پاداش، دو State value و Action Value به‌عنوان معیاری از مناسب بودن State ها و Action ها به‌صورت زیر تعریف می‌شوند:

- **State Value**: مقدار پاداشی است که انتظار می‌رود با شروع از وضعیت S و دنبال کردن سیاست π به آن دست‌یابیم.

$$V_{\pi}(s) \doteq E_{\pi}[G_t | S_t = s] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right]$$

¹ Action

² Reward

³ Discount Factor

⁴ Greedy

⁵ Hyper Parameter

- **Action Value**: مقدار پاداشی است که انتظار می‌رود با شروع از وضعیت s ، انجام عمل A و دنبال کردن سیاست π به آن دست‌یابیم.

$$Q_{\pi}(s, a) \doteq E_{\pi}[G_t | S_t = s, A_t = a] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right]$$

۵-۳-۱-۴- معادلات بلمن^۱

در فرایند تصمیم‌گیری مارکوف، معادلات بلمن، **State Value** و **Action Value** ها را در زمان‌های مختلف به‌طور بازگشتی به هم مرتبط می‌سازد. برای هر یک از **State Value** ها می‌توان نوشت:

$$\begin{aligned} V_{\pi}(s) &\doteq E_{\pi}[G_t | S_t = s] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s \right] \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) [r + \gamma E_{\pi}[G_{t+1} | S_{t+1} = s']] \\ &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) [r + \gamma V_{\pi}(s')] \end{aligned}$$

همچنین، برای هر یک از **Action Value** ها داریم:

$$\begin{aligned} Q_{\pi}(s, a) &\doteq E_{\pi}[G_t | S_t = s, A_t = a] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \right] \\ &= \sum_{s'} \sum_r p(s', r|s, a) \left[r + \gamma \sum_{a'} \pi(a'|s') E_{\pi}[G_{t+1} | S_{t+1} = s', A_{t+1} = a'] \right] \\ &= \sum_{s'} \sum_r p(s', r|s, a) \left[r + \gamma \sum_{a'} \pi(a'|s') Q_{\pi}(s', a') \right] \end{aligned}$$

به این ترتیب، به دو رابطه‌ی بازگشتی می‌رسیم که **State Value** و **Action Value** ها را در زمان‌های مختلف به هم مرتبط می‌سازد. به این معادلات، معادلات بلمن می‌گویند. بار دیگر، این دو معادله را بازنویسی می‌کنیم:

$$\begin{aligned} V_{\pi}(s) &= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r|s, a) [r + \gamma V_{\pi}(s')] \\ Q_{\pi}(s, a) &= \sum_{s'} \sum_r p(s', r|s, a) \left[r + \gamma \sum_{a'} \pi(a'|s') Q_{\pi}(s', a') \right] \end{aligned}$$

¹ Bellman Equations

۵-۳-۲- شیوهی پیاده‌سازی تابع سیاست^۱ در محیط‌های گسسته

پیش‌تر بیان شد که سیاست، استراتژی است که Agent آن را در محیط دنبال می‌کند و با توجه به آن در State های مختلف، Action های مناسب را انتخاب می‌نماید. برای محیط‌های گسسته‌ی ساده با تعداد State و Action محدود، تابع هدف به‌صورت یک جدول قابل پیاده‌سازی است که سرهای آن State ها و ستون‌هایش Action های متناظر هر State می‌باشد. به این روش پیاده‌سازی تابع سیاست روش جدولی^۲ می‌گویند. Agent برای تصمیم‌گیری در هر State می‌تواند به ردیف متناظر State در جدول مراجعه نماید و با توجه به Action Value های موجود عمل کند.

باوجود سادگی پیاده‌سازی، روش جدولی برای محیط‌های پیچیده با تعداد زیاد State و Action عملی نیست زیرا جدول متناظر بسیار بزرگ خواهد بود و ذخیره‌سازی آن دشوار است. در این موارد به جای استفاده از جدول، از تخمین زن‌های تابعی^۳ استفاده می‌شود که در ورودی خود، State را دریافت می‌کنند و در خروجی Action Value های متناظر جهت تصمیم‌گیری را می‌دهند. استفاده از شبکه‌های عصبی مصنوعی^۴ به‌عنوان قوی‌ترین تخمین زن‌های تابعی، در یادگیری تقویتی جهت ذخیره‌ی سیاست بسیار متداول و مرسوم است. این شبکه‌ها، به علت توانایی بی‌نظیرشان در ایجاد نگاشت‌های غیرخطی پیچیده و سهولت فرایند به‌روزرسانی پارامترهایشان در سال‌های اخیر به یکی از اجزاء جدایی‌ناپذیر یادگیری تقویتی بدل شده‌اند.

۵-۳-۳- یادگیری در محیط‌های گسسته

در یادگیری تقویتی، فرایند یادگیری^۵ به صورت فراگرفتن رفتار صحیح در محیط در نتیجه‌ی تعامل با آن تعریف می‌شود. در فرایندهای یادگیری تقویتی گسسته، فرایند یادگیری به دو بخش کلی تقسیم می‌شود:

بخش ۱: فرایند ارزیابی سیاست^۶

با مشخص بودن سیاست π و به کمک معادلات فوق، از روش‌های گوناگونی از جمله برنامه‌نویسی پویا یا روش مونت کارلو می‌توان مقادیر $State Value$ و $Action Value$ ها را تعیین کرد. به این فرایند به‌اصطلاح ارزیابی سیاست می‌گویند. با به دست آوردن مقادیر $State Value$ و $Action Value$ ها، Agent می‌تواند در هر مرحله، بهترین عمل متناسب با یک سیاست π را انجام دهد.

بخش ۲: بهبود سیاست^۷

به فرایند ارتقاء سیاست برای رسیدن به سیاستی بهتر، اطلاق می‌شود.

¹ Policy Function

² Tabular Method

³ Function Approximator

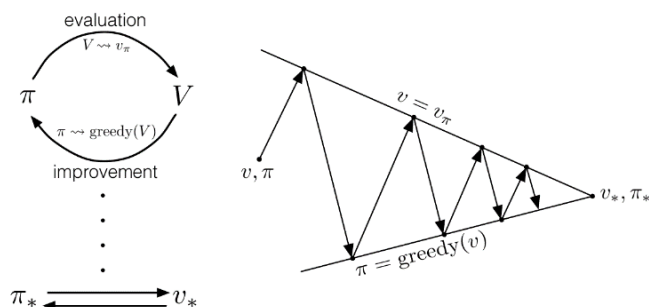
⁴ Artificial Neural Networks

⁵ Learning

⁶ Policy Evaluation

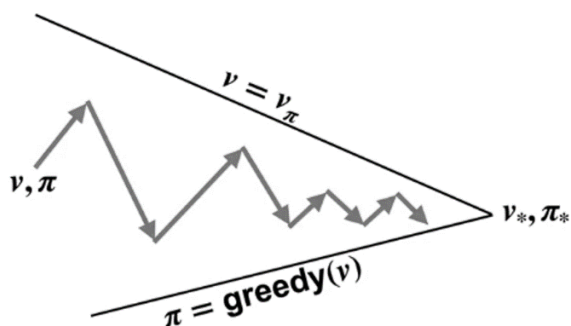
⁷ Policy Improvement

در الگوریتم‌های پایه‌ای یادگیری تقویتی، این دو فرایند مستقلاً از یکدیگر دنبال می‌شوند و بعد از هر مرحله از ارزیابی سیاست، متناسب با مقادیر به‌دست‌آمده برای Action Value و State Value ها سیاست بهبود پیدا می‌کند. این روند به‌صورت تکراری^۱ انجام می‌شود تا به یک سیاست قابل‌قبول برسیم. شکل زیر نمایی از این فرایند را نشان می‌دهد:



شکل ۴۴: فرایند تکراری ارزیابی و بهبود سیاست^۲

در مواجهه با مسائل بزرگ با محیط‌های پیچیده این روند عملی نیست زیرا تعداد State ها و Action ها معمولاً بسیار زیاد است و به‌دست آوردن مقادیر Action Value و State Value فرایندی بسیار طولانی است. در عوض، در این محیط‌ها این فرایندها به‌طور هم‌زمان و در کنار هم انجام می‌شوند تا به‌تدریج به سیاست بهینه در محیط نزدیک شویم. شکل زیر نمایانگر این فرایند است:



شکل ۴۵: اجرای هم‌زمان ارزیابی و بهبود سیاست^۳

الگوریتم‌هایی که به شکل ثانویه عمل می‌کنند، عملاً امکان یادگیری در محیط‌های پیچیده‌ی گسسته را فراهم می‌نمایند. در اینجا، ما به یکی از مهم‌ترین این الگوریتم‌ها به نام Q-Learning اشاره می‌کنیم.

^۱ Iterative

^۲ برگرفته از coursera.org

^۳ برگرفته از coursera.org

۵-۳-۱- الگوریتم Q-Learning

همان‌گونه که از نامش مشخص است، در Q-Learning به یادگیری Action Value ها پرداخته می‌شود.

فرمولاسیون به‌روزرسانی مقادیر Action Value در Q-Learning به شرح زیر است:

$$Q^{new}(s_t, a_t) = (1 - \alpha)Q^{old}(s_t, a_t) + \alpha(r_t + \gamma \max_a Q(s_{t+1}, a))$$

در عبارت فوق:

- γ همان نرخ کاهش^۱ است.

- α یک عدد مثبت بین صفر و یک موسوم به ضریب یادگیری^۲ است.

به ازای ضرایب یادگیری بزرگ ($\alpha \rightarrow 1$)، *Agent* بدون توجه به دانش کسب‌شده‌ی پیشین خود، مقدار *Action Value* را تنها بر اساس مشاهدات حال حاضرش به‌روزرسانی می‌نماید. در نقطه‌ی مقابل، به ازای ضرایب یادگیری کوچک ($\alpha \rightarrow 0$)، مقادیر *Action Value* تغییری نمی‌کنند و یادگیری انجام نمی‌شود.

به‌این ترتیب، شبه کد^۳ الگوریتم Q-Learning با استفاده از سیاست *Epsilon-greedy* به‌صورت زیر است:

شبه کد ۱: الگوریتم Q-Learning

Require:

Sates $\mathcal{X} = \{1, \dots, n_x\}$

Actions $\mathcal{A} = \{1, \dots, n_a\}$, $A : \mathcal{X} \Rightarrow \mathcal{A}$

Reward function $R : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$

Black-box (probabilistic) transition function $T : \mathcal{X} \times \mathcal{A} \rightarrow \mathcal{X}$

Learning rate $\alpha \in [0, 1]$, typically $\alpha = 0.1$

Discounting factor $\gamma \in [0, 1]$

procedure QLEARNING($\mathcal{X}, A, R, T, \alpha, \gamma$)

Initialize $Q : \mathcal{X} \times \mathcal{A} \rightarrow \mathbb{R}$ arbitrarily

while Q is not converged **do**

Start in state $s \in \mathcal{X}$

while s is not terminal **do**

Calculate π according to Q and exploration strategy (e.g. $\pi(x) \leftarrow$

$\arg \max_a Q(x, a)$)

$a \leftarrow \pi(s)$

$r \leftarrow R(s, a)$

▷ Receive the reward

$s' \leftarrow T(s, a)$

▷ Receive the new state

$Q(s', a) \leftarrow (1 - \alpha) \cdot Q(s, a) + \alpha \cdot (r + \gamma \cdot \max_{a'} Q(s', a'))$

$s \leftarrow s'$
return Q

در *Q-Learning* معمولی از روش‌های جدولی^۴ جهت پیاده‌سازی سیاست استفاده می‌شود.

¹ Discount Factor

² Learning Rate

³ Pseudocode

⁴ Tabular Methods

۵-۳-۲- الگوریتم Deep Q-Learning

Deep Q-Learning استفاده از یک شبکه‌ی عصبی عمیق^۱ جهت پیاده‌سازی سیاست در کنار قانون به‌روزرسانی مطرح‌شده در Q-Learning است. به این شبکه اصطلاحاً Policy Network گویند. ورودی این شبکه، یک State به‌خصوص و خروجی آن Action Value های متناظر Action های آن State می‌باشد. در این الگوریتم، به‌روزرسانی مقادیر Action Value عملاً معادل به‌روزرسانی وزن‌های شبکه عصبی است.

تابع هزینه‌ی تعریف‌شده برای آموزش شبکه عصبی در این الگوریتم به فرم زیر است:

$$Loss = \left[R_{t+1} + \gamma \max_{a'} q(s', a') \right] - q(s, a)$$

این تابع، معادل قانون به‌روزرسانی برای Q-Learning معمولی است. برای محاسبه‌ی تابع هزینه‌ی فوق، می‌توان مراحل زیر را دنبال می‌کنیم:

- گام ۱: مقدار Action Value های یک State را به کمک شبکه‌ی عصبی محاسبه می‌نماییم.
- گام ۲: بسته به نوع سیاست (حریصانه، epsilon-greedy یا ...) Action موردنظر را انتخاب می‌نماییم.
- گام ۳: Action را انجام می‌دهیم و به State بعدی می‌رویم (s'). همچنین، پاداش R_{t+1} را از محیط دریافت می‌کنیم.
- گام ۴: مقادیر Action Value برای State جدید را به کمک شبکه محاسبه می‌نماییم. و بیشینه‌ی آن ($\gamma \max_{a'} q(s', a')$) را به‌دست می‌آوریم.
- گام ۵: Loss را حساب می‌کنیم.

روش فوق، یک ایراد اساسی دارد: وزن‌های شبکه هم با توجه به $q(s, a)$ و هم با توجه به $q(s', a')$ به‌روزرسانی می‌شوند. این موضوع در کار شبکه اختلال ایجاد می‌کند زیرا مثلاً برای حالتی که $R_{t+1} + \gamma \max_{a'} q(s', a')$ بزرگتر از $q(s, a)$ باشد، وزن‌ها طوری به‌روزرسانی می‌شوند که $\max_{a'} q(s', a')$ کوچک و $q(s, a)$ بزرگ شود. این در حالی است که ما فقط می‌خواهیم به‌روزرسانی شبکه در راستای بزرگ شدن $q(s, a)$ باشد. جهت رفع این مشکل، یک کپی از شبکه‌ی اصلی تعیین می‌کنیم، وزن‌های آن را ثابت نگه می‌داریم و مقدار $\max_{a'} q(s', a')$ را به کمک آن حساب می‌کنیم. به این شبکه اصطلاحاً Target Network می‌گویند. بعد از گذشت تعداد مشخصی اپیزود، مقدار Target Network را دوباره برابر Policy Network قرار می‌دهیم.

در هنگام آموزش شبکه‌های عصبی عمیق، معمولاً چند داده به‌طور هم‌زمان وارد شبکه می‌شوند، گرادیان روی تمامی داده‌ها حساب می‌شود و پارامترهای شبکه به کمک میانگینی از گرادیان به‌دست‌آمده به‌روزرسانی می‌شود. این کار دو مزیت بزرگ دارد: اول آنکه استفاده از میانگینی از چند گرادیان به جای یک گرادیان باعث کاهش نویز گرادیان شده و همگرایی را تسریع می‌نماید. همچنین، این کار باعث جلوگیری از تقلید رفتار یک داده‌ی خاص توسط شبکه می‌شود. آموزش شبکه به کمک چند داده و به‌صورت هم‌زمان موجب استفاده از پردازش

¹ Deep Neural Network

موازی، بهره‌برداری مناسب از توان محاسباتی و کوتاه شدن فرایند آموزش می‌شود. برای استفاده از این خاصیت در آموزش شبکه‌های عصبی مورد استفاده در Q-Learning، مفهومی تحت عنوان Replay Memory برای الگوریتم تعریف می‌شود. Replay Memory عبارت است از یک لیست از وضعیت‌های مشاهده‌شده قبلی و پاداش دریافتی از محیط که در هر زمان از به‌روزرسانی شبکه تعدادی از آن‌ها ذخیره و استفاده می‌شود. به عبارت دیگر، در Replay Memory تعدادی لیست به فرم زیر ذخیره می‌شود:

$$(s_t, a_t, r_{t+1}, s_{t+1})$$

در هر بار تعامل با محیط، یک لیست به فرم زیر به Replay Memory اضافه می‌شود. همچنین، در هر مرحله از به‌روزرسانی وزن‌ها، تعدادی از لیست‌های فوق انتخاب‌شده و یک دسته داده وارد شبکه می‌شوند. فرایند به‌روزرسانی وزن‌های شبکه روی این دست داده انجام می‌شود.

به این ترتیب، شبه کد^۱ الگوریتم Deep Q-Learning با استفاده از سیاست *Epsilon-greedy* به صورت زیر است:

شبه کد ۲: الگوریتم Deep Q-Learning

```

Initialize replay memory  $\mathcal{D}$  to capacity  $N$ 
Initialize action-value function  $Q$  with random weights
for episode = 1,  $M$  do
    Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$ 
    for  $t = 1, T$  do
        With probability  $\epsilon$  select a random action  $a_t$ 
        otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$ 
        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$ 
        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$ 
        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$ 
        Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$ 
        Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$ 
        Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to equation 3
    end for
end for

```

۴-۵- مسائل یادگیری تقویتی پیوسته

در بسیاری از مسائل یادگیری تقویتی، Agent باید Action خود را از یک فضای پیوسته انتخاب کند. در چنین مسائلی، تغییر مقادیر Action حتی به مقدار اندک ممکن است موجب تغییرات بزرگ در محیط و شکست Agent در مأموریتش شود پس

¹ Pseudocode

ضروری است تا بتوان فضای Action را به صورت پیوسته مدل سازی نمود. روش های گسسته ای که تا به اینجای کار مورد بررسی قرار گرفت، از جمله Deep Q-Learning، دارای Action های گسسته هستند و بنابراین جهت حل مسائل یادگیری تقویتی پیوسته مناسب نیستند.

زمانی که فضای Action گسسته است، می توان با محاسبه ی Action Value برای هر Action و یافتن Action با بیشترین Action Value تصمیمات بهینه در راستای سیاست گرفت. برای پیاده سازی رویکردی مشابه در فضای پیوسته، لازم است تا یک تابع پیوسته از Action ها و State ها مانند $Q(s, a)$ داشته باشیم و در هنگام قرارگیری در یک State مانند s_0 ، مقداری از Action را انتخاب نماییم که تابع $Q(s_0, a)$ را بیشینه نماید. این به معنی حل یک مسئله ی بهینه سازی در هر بار فرایند تصمیم گیری است. واضح است که این روش برای محیط های پیوسته به علت طولانی و پیچیده بودن محاسبات به هیچ عنوان عملی نیست.

۵-۴-۱- فرایند یادگیری در محیط پیوسته

از مهم ترین الگوریتم های یادگیری تقویتی پیوسته می توان به الگوریتم Deep Deterministic Policy Gradient (DDPG) اشاره نمود. در ادامه به توضیح این الگوریتم می پردازیم.

۵-۴-۱-۱- الگوریتم Deep Deterministic Policy Gradient (DDPG)

الگوریتم Deep Deterministic Policy Gradient (DDPG) را می توان به عنوان همتای پیوسته ی Deep Q-Learning در نظر گرفت. به مانند Q-Learning، این الگوریتم بر پایه ی این نکته استوار است که اگر مقدار $Q^{\pi^*}(s, a)$ را داشته باشیم، برای هر State بهترین Action به صورت زیر قابل تعیین است:

$$a^{\pi^*}(s) = \operatorname{argmax}_a Q^{\pi^*}(s, a)$$

بنابراین الگوریتم DDPG سعی به یافتن Action بهینه برای هر وضعیت از طریق یافتن $Q^{\pi^*}(s, a)$ دارد.

این الگوریتم به طور هم زمان دو تابع را فرامی گیرد:

۱- یک تابع به نام Q-function به فرم $Q(s, a)$ که ورودی آن یک State مانند s و یک Action مانند a و

خروجی آن Action Value متناظر است.

۲- یک تابع سیاست به فرم $\mu(s)$ که ورودی آن یک State و خروجی آن Action منتخب است.

زمانی که فضای Action گسسته است، می توان با محاسبه ی Action Value برای هر Action و یافتن Action با بیشترین Action Value تصمیمات بهینه در راستای سیاست گرفت. برای پیاده سازی رویکردی مشابه در فضای پیوسته، لازم است تا یک تابع پیوسته از Action ها و State ها مانند $Q(s, a)$ داشته باشیم و در هنگام قرارگیری در یک State مانند s_0 ، مقداری از Action را انتخاب نماییم که تابع $Q(s_0, a)$ را بیشینه نماید. این به معنی حل یک

مسئله‌ی بهینه‌سازی در هر بار فرایند تصمیم‌گیری است. واضح است که این روش برای محیط‌های پیوسته به علت طولانی و پیچیده بودن محاسبات به‌هیچ‌عنوان عملی نیست.

از آنجاکه فضای Action پیوسته است، تابع $Q^*(s, a)$ نسبت به a مشتق‌پذیر است. این زمینه‌ی اجرای یک الگوریتم مبتنی بر گرادینان مؤثر جهت محاسبه‌ی یک سیاست به فرم $\mu(s)$ را می‌دهد. این تابع، به‌عنوان ورودی، یک State می‌گیرد و در خروجی، Action انتخابی متناظر آن State و مطابق با سیاست را می‌دهد. در ادامه، به جای اجرای یک فرایند بهینه‌سازی جهت یافتن مقدار a که تابع $Q(s_0, a)$ را بیشینه می‌کند، می‌توان بیشینه‌ی این تابع را با مقدار $Q(s_0, \mu(s))$ تخمین زد.

جهت بررسی دقیق‌تر الگوریتم DDPG، ابتدا معادله‌ی بهینه‌ی بلمن برای Action Value ها را در نظر می‌گیریم:

$$Q^{\pi^*}(s, a) = E_{s' \sim P}[r(s, a) + \gamma \max_a Q^{\pi^*}(s', a')]$$

در عبارت فوق، منظور از $s' \sim P$ آن است که State بعدی یا s' از یک توزیع احتمالاتی به فرم $P(\cdot | s, a)$ می‌آید. معادله‌ی فوق نقطه‌ی شروعی برای یافتن یک تخمین‌زن برای $Q^{\pi^*}(s, a)$ است. فرض کنید تخمین‌زن به فرم $Q_\phi(s, a)$ در دست داریم که ϕ پارامترهای تخمین‌زن است. تابع هزینه‌ی^۱ مورد استفاده جهت سنجش عملکرد تخمین‌زن به فرم زیر می‌باشد:

$$L = \left(Q_\phi(s, a) - \left(r + \gamma \max_{a'} Q_\phi(s', a') \right) \right)^2$$

چالش اصلی در محاسبه‌ی عبارت فوق، محاسبه‌ی مقدار $\max_{a'} Q_\phi(s', a')$ است. در الگوریتم DDPG فرض می‌شود:

$$\max_{a'} Q_\phi(s', a') = \mu(s')$$

و بنابراین خواهیم داشت:

$$L = \left(Q_\phi(s, a) - (r + \gamma \mu(s')) \right)^2$$

ذکر دو نکته‌ی پایانی در مورد این الگوریتم خالی از لطف نیست:

۱- بنا به دلایلی کاملاً مشابه با آنچه برای Deep Q-Learning گفته شد، استفاده از Replay Memory در الگوریتم DDPG نیز مرسوم و رایج است.

۲- به‌مانند Deep Q-Learning، در DDPG نیز از یک Target Network جهت به‌روزرسانی $Q_\phi(s, a)$ استفاده می‌شود.

¹ Loss Function

در الگوریتم DDPG، گاهی اوقات به شبکه‌ای که وظیفه‌ی تعیین Action از State ورودی را دارد، Actor و به شبکه‌ای که وظیفه‌ی تعیین Action Value متناظر را دارد Critic می‌گویند.

شبه کد الگوریتم DDPG در ادامه آورده شده است:

شبه کد ۳: الگوریتم Deep Deterministic Policy Gradient (DDPG)

```

1: Input: initial policy parameters  $\theta$ , Q-function parameters  $\phi$ , empty replay buffer  $\mathcal{D}$ 
2: Set target parameters equal to main parameters  $\theta_{\text{targ}} \leftarrow \theta$ ,  $\phi_{\text{targ}} \leftarrow \phi$ 
3: repeat
4:   Observe state  $s$  and select action  $a = \text{clip}(\mu_{\theta}(s) + \epsilon, a_{\text{Low}}, a_{\text{High}})$ , where  $\epsilon \sim \mathcal{N}$ 
5:   Execute  $a$  in the environment
6:   Observe next state  $s'$ , reward  $r$ , and done signal  $d$  to indicate whether  $s'$  is terminal
7:   Store  $(s, a, r, s', d)$  in replay buffer  $\mathcal{D}$ 
8:   If  $s'$  is terminal, reset environment state.
9:   if it's time to update then
10:    for however many updates do
11:      Randomly sample a batch of transitions,  $B = \{(s, a, r, s', d)\}$  from  $\mathcal{D}$ 
12:      Compute targets

$$y(r, s', d) = r + \gamma(1 - d)Q_{\phi_{\text{targ}}}(s', \mu_{\theta_{\text{targ}}}(s'))$$

13:      Update Q-function by one step of gradient descent using

$$\nabla_{\phi} \frac{1}{|B|} \sum_{(s, a, r, s', d) \in B} (Q_{\phi}(s, a) - y(r, s', d))^2$$

14:      Update policy by one step of gradient ascent using

$$\nabla_{\theta} \frac{1}{|B|} \sum_{s \in B} Q_{\phi}(s, \mu_{\theta}(s))$$

15:      Update target networks with

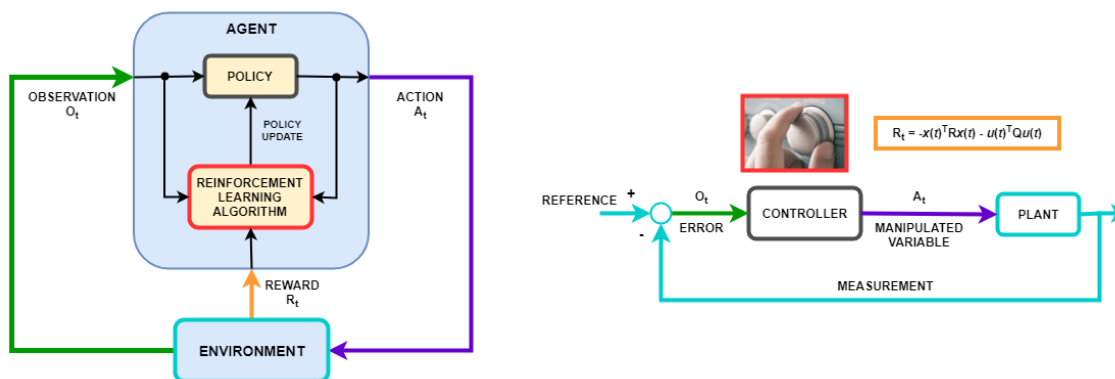
$$\begin{aligned} \phi_{\text{targ}} &\leftarrow \rho \phi_{\text{targ}} + (1 - \rho) \phi \\ \theta_{\text{targ}} &\leftarrow \rho \theta_{\text{targ}} + (1 - \rho) \theta \end{aligned}$$

16:    end for
17:  end if
18: until convergence

```

۵-۵- کاربرد یادگیری تقویتی در مسائل حوزه مهندسی کنترل

رفتار یک تابع سیاست در یادگیری تقویتی (مشاهده‌ی وضعیت محیط و تصمیم‌گیری برای انجام اعمال مختلف بر اساس آن) شباهت به نقش کنترلر در مسائل حوزه‌ی کنترل دارد. یادگیری تقویتی می‌تواند به صورت یک مسئله‌ی کنترل در نظر گرفته شود. شکل زیر نمایانگر نزدیکی یادگیری تقویتی به فرایند کنترل با بازخورد^۱ است.



شکل ۴۶: تشابه یادگیری تقویتی و فرایند کنترل با بازخورد^۲

ارکان اساسی مسائل این دو حوزه، مطابق جدول زیر به هم ارتباط دارند:

جدول ۱۲: تشابه مفاهیم حوزه‌ی کنترل و یادگیری تقویتی

یادگیری تقویتی	کنترل سیستم‌ها
سیاست (Policy)	کنترلر (Controller)
محیط (Environment)	هر آنچه کنترلر نیست (شامل Plant، فیلتر، نویز اندازه‌گیری، مبدل آنالوگ به دیجیتال و برعکس و ...)
مشاهده (Observation)	اندازه‌گیری (Measurement)
عمل (Action)	متغیر کنترلی (Control Variable)
پاداش (Reward)	سیگنال خطای ورودی به کنترلر (Error Signal)

به این ترتیب، یادگیری تقویتی را می‌توان به عنوان یک کنترلر در محیطش در نظر گرفت. برخلاف کنترلرهای کلاسیک (مانند PID) برای طراحی کنترلر مبتنی بر یادگیری تقویتی هیچ نیازی به دانستن دینامیک محیط نیست و مستقل از پیچیدگی معادلات حاکم بر سیستم، می‌توان از یادگیری تقویتی جهت کنترل آن استفاده نمود.

^۱ Feedback

^۲ برگرفته از Mathworks.com

۶- کنترل آونگ معکوس دورانی به کمک یادگیری تقویتی

چندین مسئله‌ی کنترلی را می‌توان برای سازه‌ی آونگ دورانی معکوس متصور شد. در ساده‌ترین حالت، صرفاً موقعیت زاویه‌ای آونگ حایز اهمیت و به موقعیت و نوع حرکت بازو توجهی نمی‌شود. این مسئله، یک مسئله‌ی کنترلی تک ورودی - تک خروجی^۱ است. در مسائل کنترلی پیچیده‌تر، می‌توان موقعیت بازو را نیز وارد سیستم کرد تا مسئله‌ی کنترلی به یک مسئله‌ی چند ورودی - تک خروجی^۲ بدل شود. در این حالات، می‌توان موقعیت زاویه‌ای دلخواهی برای بازو تعیین نمود تا دقیقاً در آن قرار گیرد. همچنین، می‌توان حرکت بازو را طوری محدود کرد تا با حداقل حرکت، آونگ را کنترل نماید.

در این پروژه، به علت ضیق وقت، تنها به بررسی عملکرد یادگیری تقویتی روی مسئله‌ی اول (کنترل وضعیت آونگ بدون توجه به وضعیت بازو) می‌پردازیم و بررسی حالات دیگر مسئله‌ی کنترل آونگ دورانی معکوس را به پژوهش‌های بعدی موکول می‌نماییم. برای کنترل آونگ معکوس دورانی، به کمک یادگیری تقویتی ابتدا لازم است تا ورودی‌ها و خروجی‌های سیستم را مشخص نماییم:

- ورودی سیستم کنترلی دوره‌ی کاری^۳ سیگنال PWM ورودی به سیستم و جهت چرخش موتور می‌باشد.

- خروجی سیستم کنترلی موقعیت آونگ در صفحه‌ی چرخش می‌باشد.

جهت تعریف یک مسئله‌ی یادگیری تقویتی برای کنترل آونگ معکوس دورانی با دو مسئله‌ی زیر مواجه هستیم:

۱- تعیین ورودی الگوریتم کنترلی

۲- تعیین خروجی الگوریتم کنترلی

در ادامه، در دو بخش مجزا سعی می‌شود به مسائل فوق پاسخ مناسب داده شود.

۶-۱- ورودی الگوریتم کنترلی

یک رویکرد ممکن برای ورودی سیستم کنترلی این است که در هر لحظه تنها موقعیت آونگ به الگوریتم داده شود. اما با این رویکرد هرگز نمی‌توان کنترلر پایدار مدنظرمان را بسازیم. علت این است که آونگ معکوس دورانی، یک سیستم با مرتبه‌ی مخالف صفر است و وضعیت آونگ در هر زمان به وضعیت گذشته‌اش وابسته است. بنابراین برای نمایش کامل وضعیت آونگ، لازم است تا علاوه بر وضعیت زمان حال، اطلاعاتی در مورد گذشته‌ی آن نیز در اختیار الگوریتم قرار داده شود. این کار می‌تواند به دو صورت انجام شود:

۱- موقعیت زاویه‌ای آونگ و تعداد دلخواهی از مشتقات عددی آن (سرعت زاویه‌ای، شتاب زاویه‌ای و ...) به صورت یک بردار به شبکه داده شود.

۲- یک سری زمانی به طول دلخواه از موقعیت‌های اخیر آونگ به شبکه داده شود.

¹ Single Input, Single Output (SISO)

² Multiple Inputs, Single Output (MISO)

³ Duty Cycle

استفاده از روش اول به نظر منطقی‌تر است زیرا کاری که در این روش انجام می‌پذیرد باعث مستقل شدن ورودی‌ها از زمان نمونه‌برداری^۱ می‌شود. این درحالی است که در الگوریتمی با ورودی‌های از نوع دوم، زمان نمونه‌برداری نمی‌تواند تغییر کند و به عنوان مثال، الگوریتمی که با داده‌های نمونه‌برداری شده با فرکانس 1000 Hz آموزش داده شده است نمی‌تواند خروجی صحیحی به ازای ورودی‌هایی از فرکانس‌های نمونه برداری دیگر تولید کند. علت این موضوع آن است که فاصله‌ی زمانی بین داده‌های موجود در سری زمانی حائز اهمیت است. بنا به دلیل ذکر شده، از روش اول برای تولید ورودی‌های الگوریتم استفاده می‌شود.

۶-۲- خروجی الگوریتم کنترلی

بدیهی است که ورودی الگوریتم یادگیری تقویتی مورد استفاده جهت کنترل سیستم (دوره‌ی کاری سیگنال PWM ورودی به سیستم)، مقداری پیوسته است. به این ترتیب، جهت حل طراحی یک کنترلر مناسب برای آونگ معکوس دورانی به کمک یادگیری تقویتی، دو راهکار کلی وجود دارد:

- **راهکار اول:** فضای تصمیم‌گیری (مقادیر دوره‌ی کاری سیگنال PWM ورودی به سیستم) را گسسته‌سازی^۲

کنیم. جهت دریافت جواب مناسب از سیستم در این حالت، لازم است تا تعداد مقادیر گسسته نهایی به اندازه‌ی کافی زیاد باشد تا بتواند به درستی تخمینی از حالت پیوسته را ارائه دهد. در این راهکار، می‌توان مسئله را با الگوریتم Deep Q-Learning حل نماییم.

- **راهکار ثانویه:** بدون گسسته‌سازی فضای تصمیم‌گیری، مستقیماً از روش‌های یادگیری تقویتی پیوسته جهت

کنترل آونگ معکوس دورانی استفاده نماییم. در این حالت می‌توانیم مسئله را به کمک الگوریتم Deep Deterministic Policy Gradient (DDPG) حل نماییم.

هر یک از دو راهکار فوق مزایا و معایب مختص به خود را دارند و تا قبل از مشاهده‌ی نتایج عملی عملکردشان نمی‌توان به طور قطعی نتیجه‌گیری در مورد آنها انجام داد.

- پیاده‌سازی الگوریتم Deep Q-Learning و آموزش آن به مراتب آسان‌تر از پیاده‌سازی الگوریتم DDPG و فرایند

آموزش آن است. علت این موضوع، ساختار پیچیده‌تر الگوریتم DDPG می‌باشد و همانطور که پیش‌تر توضیح داده شده بود، در DDPG لازم است تا دو شبکه جهت دریافت Action انتخابی متناسب با State و Value متناظر با State و Action به صورت مجزا پیاده‌سازی و آموزش داده شوند. این درحالی است که در Q-Learning تنها یک شبکه Policy Network آموزش داده می‌شود.

- با وجود آسانی پیاده‌سازی Deep Q-Learning، لازم است تا جهت رسیدن به تقریب مناسبی از محیط Action

پیوسته، تعداد Action‌های گسسته بیشتر شود. این موضوع، موجب می‌شد که ساختار شبکه‌ی عصبی در Q-Learning اندکی پیچیده‌تر از ساختار شبکه‌ی عصبی در DDPG شود. با توجه به این که با یک مسئله‌ی کنترل بی‌درنگ مواجه هستیم، پیچیده‌تر شدن ساختار می‌تواند موجب ایجاد تاخیر زیاد در محاسبات و نهایتاً ایجاد

¹ Sampling Time

² Discretize

اشکال در فرایند کنترل آونگ شود. البته، مطلب ذکر شده بسته به میزان حساسیت سیستم به تغییرات کوچک در سیگنال کنترلی ورودی دارد. اگر در تجربه ثابت شود که سیستم به تغییرات جزئی در ولتاژ موتور حساسیت کمی دارد، می‌توان مقادیر گسسته‌ی ورودی را با فاصله‌ی بیشتری از هم انتخاب نمود تا فضای Action های گسسته کوچکتر شود.

در طی فرایند آموزش، عملکرد هر دو شبکه‌ی فوق را بر کارکرد آونگ معکوس دورانی بررسی می‌نماییم.

۷- فرایند آموزش

جهت پیاده‌سازی هر دو الگوریتم Deep Q-Learning و Deep Deterministic Policy Gradient از فریم‌ورک^۱ Keras از کتابخانه‌ی Tensor-Flow در زبان برنامه‌نویسی پایتون استفاده می‌شود. همچنین، در اجرای فرایند آموزش از یک جی‌پی‌یو^۲ Nvidia GTX-1050 استفاده می‌کنیم. فرایند آموزش به کمک هر دو الگوریتم، از دو بخش کلی تشکیل می‌شود:

- بخش ۱: آموزش در محیط شبیه‌سازی Simulink که به موازات فرایند ساخت دستگاه انجام شد
- بخش ۲: آموزش روی سازه‌ی فیزیکی آونگ معکوس دورانی. در طی این فرایند، از دانش کسب شده در محیط شبیه‌سازی به عنوان یک دانش اولیه^۳ برای کنترل سازه‌ی فیزیکی استفاده شد. به عبارتی، شبکه‌های عصبی مورد استفاده در آموزش سازه‌ی فیزیکی، معماری مشابه شبکه‌های آموزش داده شده در محیط شبیه‌سازی داشتند و با وزن‌های مشابه، وزن‌دهی اولیه شده بودند.

در فرایند آموزش به کمک هر دو الگوریتم ذکر شده در بالا، مقادیر مکان، سرعت و شتاب زاویه‌ای در هر لحظه به عنوان State انتخاب می‌شوند. همچنین به عنوان تابع پاداش، از تابعی به فرم زیر استفاده خواهیم کرد. این تابع هدف، از محیط آونگ معکوس خطی کتابخانه‌ی Gym^۴ در پایتون اخذ شده است.

$$R = - \left(\theta_*^2 + 0.1 \dot{\theta}_*^2 + 0.0001 * (motor\ voltage)^2 \right)$$

در این معادله:

- θ_* به صورت $\frac{\theta}{\pi}$ تعریف می‌شود که θ موقعیت زاویه‌ای آونگ نسبت به محور قائم است. (مطابق شکل ۱)
- $\dot{\theta}$ اندیس Action انتخابی می‌باشد.

^۱ Framework

^۲ Graphical Processing Unit (GPU)

^۳ Prior Knowledge

^۴ کتابخانه‌ی Gym یک کتابخانه‌ی گرافیکی برای Python است که در آن تعدادی از سیستم‌های فیزیکی و غیر فیزیکی جهت اجرای آموزش یادگیری تقویتی پیاده‌سازی شده‌اند.

رابطه‌ی فوق از سه بخش مجزا تشکیل شده است که به توضیح علت وجود هر سه بخش می‌پردازیم:

- ترم θ_*^2 در تابع پاداش سعی در نزدیک نگه داشتن آونگ به نقطه‌ی تعادل بالایی‌اش دارد.
- ترم $0.1\dot{\theta}_*^2$ سعی در کوچک نگه داشتن سرعت آونگ و نزدیک کردن آن به وضع پایدار دارد.
- ترم $(motor\ voltage)^2 * 0.0001$ سعی در نزدیک نگه داشتن ولتاژ پایه‌های موتور به صفر دارد.

۷-۱- الگوریتم Deep Q-Learning

همانطور که در بخش ۶ توضیح داده شده بود، جهت استفاده از الگوریتم Deep-Q Learning برای کنترل آونگ معکوس دورانی، لازم است تا فضای Action‌ها (دوره‌ی کاری سیگنال PWM ورودی به موتور) را گسسته‌سازی نماییم. بازه‌ی تغییر ولتاژ ورودی مجاز برای موتور از 24V تا -24V می‌باشد. از طرفی، مطابق آنچه در بخش ۴ گفته شد، آستانه‌ی حرکت موتور در هنگام اتصال به آونگ و بازو حدود 4.5V است. فرایند گسسته‌سازی را به صورتی انجام می‌دهیم که ۹ ورودی مجاز در این بازه داشته باشیم. به عبارتی، مقادیر گسسته برای سیگنال‌های PWM را طوری تعیین می‌کنیم که مقادیر ولتاژ متناظر موتور مطابق سری زیر شود:

$$V_1 = -24V, V_2 = -18V, V_3 = -12V, V_4 = -6V,$$

$$V_5 = 0V, V_6 = 6V, V_7 = 12V, V_8 = 18V, V_9 = 24V$$

Action‌های متناظر جهت تولید ولتاژهای بالا در پایه‌های موتور، به ترتیب a_1, a_2, a_3, \dots و a_9 است. به این ترتیب، خروجی لایه‌ی آخر شبکه‌ی عصبی استفاده شده در الگوریتم Deep-Q Learning ۹ نورون خواهد داشت.

همچنین، جهت فراگیری سیاست مناسب، از یک شبکه‌ی عصبی ۴ لایه با معماری زیر استفاده می‌شود:

Input Size: State Size (3)

Output Size: Number of Actions (9)

Architecture:

Layer 1: Dense (n_neurons=64, input_dim= Input Size, activation=Leaky ReLu)

Layer 2: Dense (n_neurons=128, input_dim= 64, activation= Leaky ReLu)

Layer 3: Dense (n_neurons=64, input_dim= 128, activation=Sigmoid)

Layer 3: Dense (n_neurons=Output Size, input_dim=32, activation=Linear)

همچنین، پارامترهای آموزش شبکه به شرح زیر می‌باشند:

- Optimizer: ADAM
- Learning Rate: $\alpha = 0.001$
- Decay Rate: $\gamma = 0.95$
- Replay Memory Size: 2000
- Initial Exploration Rate: $\epsilon = 1$
- Initial Exploration Rate: $\epsilon = 0.995$
- Minimum Exploration Rate: $\epsilon_{min} = 0.01$

پارامترهای آموزشی تعیین شده همگی از روی مدل‌های موفق برای آموزش آونگ معکوس خطی در محیط Gym الهام گرفته شده‌اند.

۷-۲- الگوریتم Deep Deterministic Policy Gradient

بر خلاف الگوریتم Deep-Q Learning، در هنگام استفاده از الگوریتم Deep Deterministic Policy Gradient دیگر نیازی به گسسته‌سازی مقادیر Action نیست. خروجی این الگوریتم، یک مقدار بین منفی یک و یک است که به کمک یک نگاشت خطی، به مقادیر متناظر برای دوره‌ی کاری سیگنال PWM تبدیل می‌شود. جهت تولید این مقادیر، در لایه‌ی آخر شبکه از یک تابع فعال‌ساز تانژانت هایپربولیک^۱ استفاده می‌شود.

معماری شبکه‌ی Actor استفاده شده در این بخش به شرح زیر می‌باشد:

Input Size: State Size (3)

Output Size: Action Size (1)

Architecture:

Layer 1: Dense (n_neurons=256, input_dim= Input Size, activation=Leaky ReLu)

Layer 2: Dense (n_neurons=512, input_dim= 256, activation= Leaky ReLu)

Layer 3: Dense (n_neurons=Output Size, input_dim=256, activation=Tanh)

معماری شبکه‌ی Critic استفاده شده اندکی پیچیده تر است. در این شبکه، ورودی State و Action به صورت مجزا وارد دو لایه از نوع Dense با نام‌های Action Dense و State Dense می‌شوند. در ادامه، خروجی این دو لایه به هم چسبانده شده و به عنوان ورودی یک شبکه‌ی عصبی دولایه‌ی دیگر داده می‌شود تا نهایتاً مقدار Action Value متناظر تولید گردد. معماری این شبکه در زیر آورده شده است:

State Dense:

Input Size: State Size (3)

Output Size: 256

Architecture:

Layer 1: Dense (n_neurons=256, input_dim= Input Size, activation=Leaky ReLu)

Layer 2: Dense (n_neurons=Output Size, input_dim= 256, activation=Leaky ReLu)

Action Dense:

Input Size: Action Size (1)

Output Size: 256

Architecture:

Layer 1: Dense (n_neurons=Output Size, input_dim= Input Size, activation=Leaky ReLu)

Trunk:

Input Size: State Dense Output Size + Action Dense Output Size

Output Size: 1

Architecture:

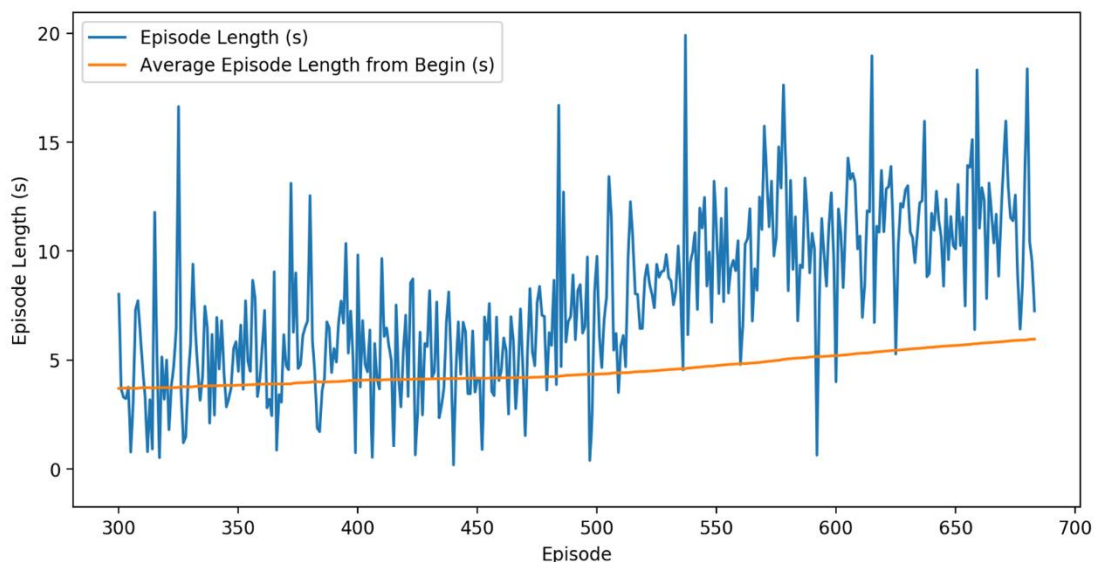
Layer 1: Dense (n_neurons=256, input_dim= Input Size, activation=Leaky ReLu)

Layer 3: Dense (n_neurons=Output Size, input_dim=256, activation=Tanh)

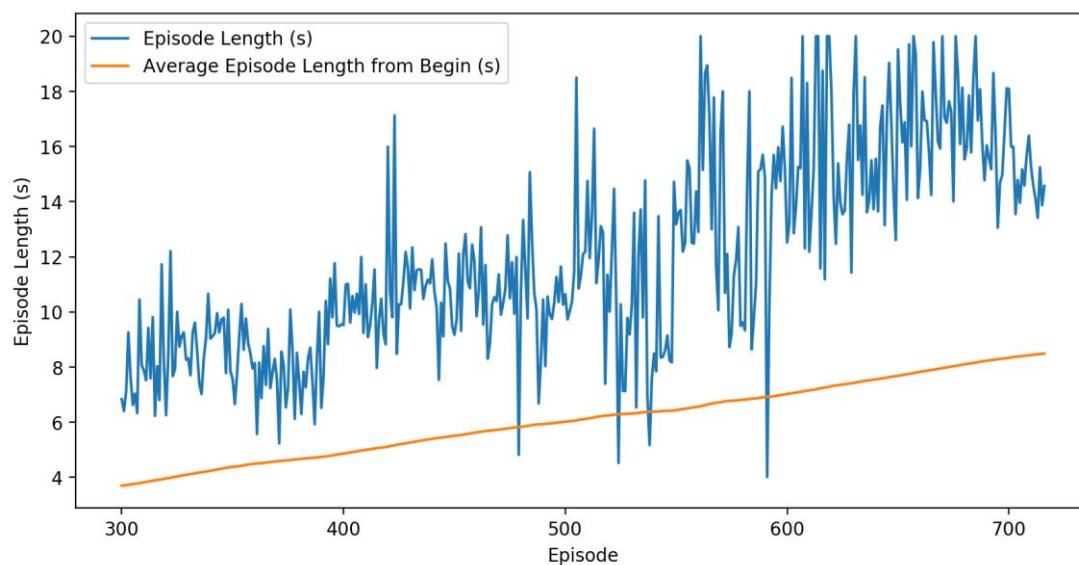
¹ Tangent Hyperbolic (Tanh)

۸- نتایج آموزش

فرایند آموزش هردو الگوریتم Deep Q-Learning و DDPG ابتدا در محیط شبیه‌سازی و برای 700 اپیاک انجام شد. نمودار مدت زمان اپیزودها برای ۴۰۰ اپیزود آخر در زیر آورده شده است:

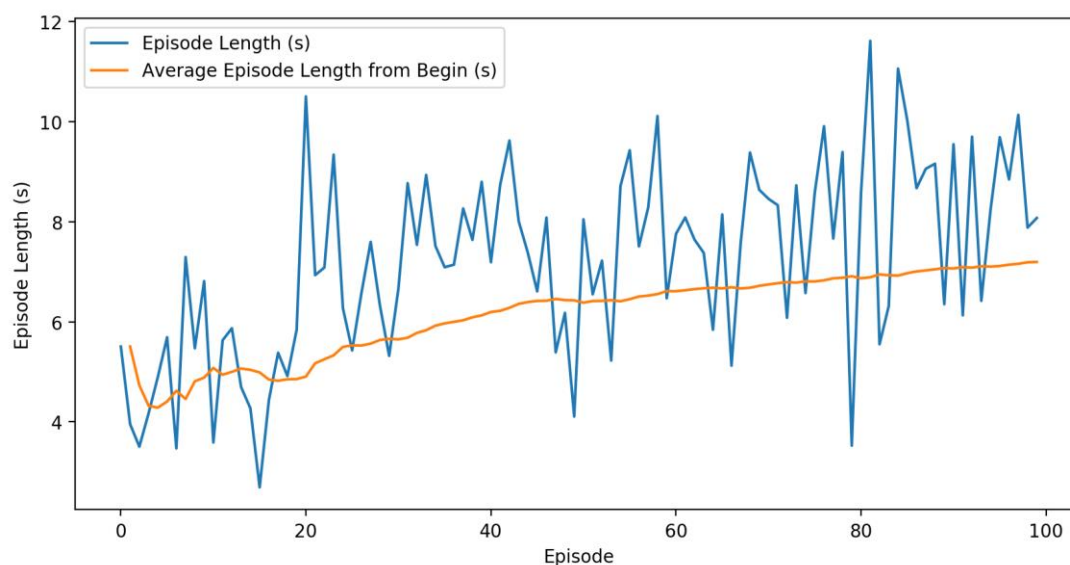


شکل ۴۷: نتایج آموزش الگوریتم *Deep Q-Learning* در ۴۰۰ اپیزود پایانی

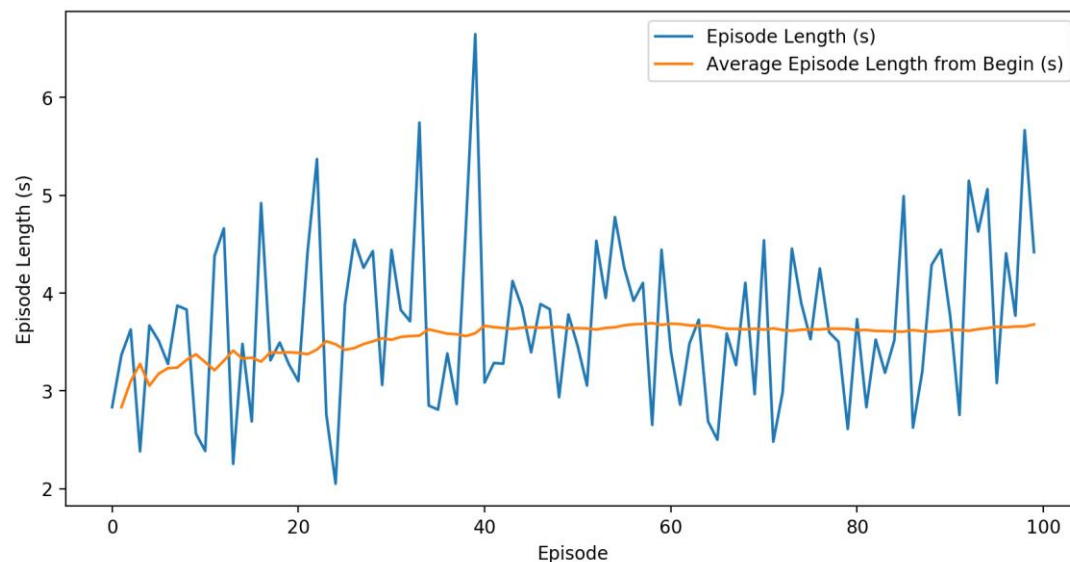


شکل ۴۸: نتایج آموزش الگوریتم *Deep Deterministic Policy Gradient* در ۴۰۰ اپیزود پایانی

در ادامه، شبکه‌های آموزش داده شده در محیط شبیه‌سازی روی سازه‌ی واقعی تست شدند. متأسفانه، به علت طولانی بودن فرایند آموزش، این فرایند در ثالب پند بخش اجرا شد و نمودار کاملی از آموزش سازه‌ی آونگ معکوس دورانی برای همه‌ی اپیزودها در دست نیست. نمودارهای زیر، عملکرد سازه‌ی آونگ برای ۱۰۰ اپیزود آموزشی برای در الگوریتم را نمایش می‌دهند.



شکل ۴۹: نتایج آموزش الگوریتم *Deep Q-Learning* در ۱۰۰ اپیزود آموزش روس سازی آونگ معکوس دورانی



شکل ۵۰: نتایج آموزش الگوریتم *Deep Deterministic Policy Gradient* در ۱۰۰ اپیزود آموزش روس سازی آونگ معکوس دورانی

۹- تفسیر و بررسی نتایج

همانطور که در بخش نتایج مشاهده شد، روش Deep Deterministic Policy Gradient در کنترل آونگ معکوس دورانی عملکرد بهتری نسبت به Deep Q-Learning از خود نشان داد. علت این موضوع را می‌توان در این نکته جستجو کرد که آونگ معکوس دورانی محیطی ذاتاً پیوسته و بسیار ناپایدار است و کوچکترین تغییر در ورودی آن می‌تواند موجب برهم خوردن تعادل و سقوطش شود.

در هنگام تست روی سازه‌ی فیزیکی، Deep Q-Learning به بازدهی بالاتری رسید. به نظر می‌رسد این موضوع نه به خاطر برتری Deep Q-Learning نسبت به Deep Deterministic Policy Gradient بلکه به علت حجم محاسبات کمتر آن باشد. قابل ذکر است که در Deep Deterministic Policy Gradient، شبکه‌ی Actor و Critic همزمان آموزش می‌بینند درحالی که در روش Deep Q-Learning تنها یک شبکه برای آموزش دادن وجود دارد. انتظار می‌رود که در صورت کوتاه‌تر کردن فرایند آموزش شبکه‌های الگوریتم Deep Deterministic Policy Gradient از طریق اصلاح معماری شبکه‌های عصبی آن و یا استفاده از پردازنده‌های قوی‌تر، این عیب الگوریتم برطرف شده و بازدهی آن در محیط واقعی نیز، مانند محیط شبیه‌سازی، از الگوریتم Deep Q-Learning پیشی گیرد.

۱۰- پیشنهادهایی جهت بهبود نتایج

در این بخش، با توجه به فعالیت‌های انجام شده در این پروژه و دانش به‌دست آمده، برخی پیشنهادهای که احتمالاً موجب بهبود عملکرد شبکه می‌شود را عنوان می‌نماییم:

۱- یکی از ارکان ضروری جهت اجرای کامل و بهینه‌ی آموزش در فرایندهای مبتنی بر یادگیری، داشتن توان محاسباتی مناسب است. کلیه‌ی فرایندهای آموزش اجرا شده در این پروژه به کمک یک جی‌پی‌یو^۱ Nvidia GTX-1050 انجام شد. همچنین، به علت لزوم برقراری سریع ارتباط با سیستم واقعی یا محیط شبیه‌ساز، امکان استفاده از سیستم‌های پردازش ابری^۲ مانند Google Colab وجود نداشت. این موضوع باعث شد تا نتوان بیشتر فرایندهای آموزش را تا نقطه‌ی بهینه ادامه داد. انتظار می‌رود با اجرای روش‌ها به کمک پردازنده‌های قوی‌تر و در طی فرایندهای طولانی‌تر آموزش، عملکرد کنترلر مبتنی بر یادگیری تقویتی بهبود زیادی داشته باشد.

۲- در این پروژه، مکان، سرعت و شتاب زاویه‌ای آونگ به عنوان State ورودی در نظر گرفته شدند. انتظار می‌رود که با وارد کردن مشتقات مرتبه بالاتر مکان زاویه‌ای آونگ در کنار مشتقات مرتبه اول و دوم به سیستم، شاهد بهبود عملکرد آن باشیم. علت این موضوع این است که با وارد کردن تعداد بیشتری از مشتقات، اطلاعات الگوریتم در مورد وضعیت پیشین سیستم، پیش‌بینی رفتار آتی آن و احتمال گرفتن تصمیم‌های صحیح بیشتر می‌شود.

¹ Graphical Processing Unit (GPU)

² Cloud Computing

۳- متاسفانه، به علت کمبود وقت، امکان بررسی تاثیر تمامی پارامترهای اساسی^۱ الگوریتم‌های Deep Q-Learning و DDPG، به مانند اندازه‌ی Replay Memory و یا مدت زمان بین بروزرسانی‌های متوالی Target Network توسط Policy Network بر عملکرد کنترلر یادگیری تقویتی میسر نشد. احتمال می‌رود که با تغییر این مقادیر بتوان به نتایج بهتر دست یافت.

۴- بهینه‌سازی معماری شبکه‌های عصبی یکی از مواردی است که انتظار می‌رود تاثیر قابل توجهی بر عملکرد هر دو الگوریتم Deep Q-Learning و DDPG داشته باشد. متاسفانه تاکنون الگوریتم مشخصی برای اصلاح معماری شبکه‌های عصبی معرفی نشده است اما استفاده از روش‌هایی مانند به تصویر کشیدن گرادینان بازگشتی و توزیع خروجی هر لایه‌ی شبکه‌ی عصبی به کمک Tensor Board می‌تواند اطلاعات مفیدی در ارتباط با تاثیر بخش‌های مختلف شبکه‌ی عصبی در اختیارمان قرار دهند. با استفاده از این داده‌ها می‌توان معماری شبکه را تا حدودی ارتقاء داد به طوری که ضمن ساده‌تر شدن، بازده عملکرد الگوریتم یادگیری تقویتی افت نکند.

۵- در مقاله‌ی ، به جای شبکه‌های عصبی Feed Forward، از لایه‌های کانولوشن^۲ تک بعدی جهت پردازش ورودی‌های شبکه استفاده شده است. عملکرد الگوریتم یادگیری تقویتی با استفاده از این معماری در محیط شبیه‌سازی، نسبت به قبل پیشرفت چشمگیری داشته است. این موضوع تا حدی قابل پیش‌بینی است زیرا شبکه‌های عصبی کانولوشنی^۳ را می‌توان به صورت شبکه‌های Feed Forward در نظر گرفت که هر نورون تنها با برخی از نورون‌های لایه‌ی قبلی ارتباط دارد. این موضوع، ضمن بدتر نکردن عملکرد شبکه نسبت به شبکه‌ی Feed Forward متناظرش، تعداد پارامترهای آن را کاهش می‌دهد که باعث می‌شود شبکه قدرت تعمیم^۴ بهتر و سرعت آموزش بیشتری داشته باشد. لذا، جایگزین کردن معماری عنوان شده برای شبکه‌ی عصبی با یک شبکه‌ی کانولوشنی به احتمال زیاد موجب بهبود فرایند کنترل آونگ می‌شود.

۶- یکی دیگر از الگوریتم‌های حوزه‌ی یادگیری تقویتی گسسته، الگوریتم Double Deep Q-Learning (DDQN) نام دارد. این الگوریتم معرفی شد تا چندی از مشکلات روش Q-Learning از جمله مشکل بیش از حد برآورد کردن Action Valueها در Deep Q-Learning و Q-Learning را برطرف سازد. به علت کمبود وقت، امکان تست عملکرد این الگوریتم میسر نشد اما پیشنهاد می‌شود جهت بهبود عملکرد روش گسسته، این روش نیز در کنار روش Deep Q-Learning امتحان گردد.

۷- متاسفانه به علت کمبود وقت، نتوانستیم عملکرد توابع پاداش^۵ مختلف را بر کارکرد الگوریتم‌ها بسنجیم. به قطع امتحان کردن توابع پاداش مختلف می‌تواند موجب بهبود عملکرد سیستم شود.

¹ Hyperparameter

² Convolution

³ Convolutional Neural Networks

⁴ Generalization

⁵ Reward

۱۱- جمع‌بندی و نتیجه‌گیری

در این پروژه به ساخت و شبیه‌سازی سیستم آونگ معکوس دورانی پرداخته شد. همچنین، عملکرد دو الگوریتم مهم حوزه‌ی یادگیری تقویتی پیوسته و گسسته (Deep Q-Learning و Deep Deterministic Policy Gradient) در شبیه‌سازی و واقعیت سنجیده شد. با وجود پیچیده‌تر و طولانی‌تر بودن فرایند آموزش Deep Deterministic Policy Gradient، این الگوریتم توانست تا حدودی از Deep Q-Learning در محیط شبیه‌سازی پیشی گیرد. علت این موضوع را می‌توان در پیوستگی ذاتی ورودی‌های آونگ معکوس دورانی و ناپایداری آن جستجو کرد. در مقابل، الگوریتم Deep Q-Learning به علت حجم کمتر محاسبات، عملکرد بهتری در کنترل سازه‌ی آونگ معکوس دورانی داشت. نتایج به‌دست آمده از الگوریتم‌های بررسی‌شده در کنترل آونگ معکوس دورانی، نشان‌دهنده‌ی توانایی این الگوریتم‌های حوزه‌ی یادگیری تقویتی در کنترل سیستم‌های غیرخطی و پیچیده است.

فعالیت‌های صورت گرفته در این پروژه، تنها گوشه‌ی کوچکی از پتانسیل‌های فراوان تحقیقاتی بر روی سازه‌ی آونگ معکوس دورانی می‌باشد. امید است فعالیت‌های صورت گرفته تاکنون، بتواند بستری مناسب برای پژوهش‌های بیشتر درمورد عملکرد الگوریتم‌های هوشمند روی آونگ معکوس دورانی، به عنوان یک مسئله‌ی معیار در حوزه‌ی مهندسی کنترل، شود.

- [1] J L Duart, B Montero, P A Ospina and E González, "*Dynamic Modeling and Simulation of a Rotational Inverted Pendulum*", Journal of Physics: Conference Series, Volume 792, VIII International Congress of Engineering Physics 7–11 November 2016, Mérida, Yucatán, Mexico.
- [2] Wudai Liao, Zhengbo Liu, Shengjun Wen, Shuhui Bi, Dongyun Wang, "*Fractional PID based stability control for a single link rotary inverted pendulum*", 2015 International Conference on Advanced Mechatronic Systems (ICAMechS).
- [3] Jia-JunWang, "*Simulation studies of inverted pendulum based on PID controllers*" Simulation Modelling Practice and Theory, Volume 19, Issue 1, January 2011, Pages 440-449.
- [4] Krishen, J., Becerra, V.M., "*Efficient fuzzy control of a rotary inverted pendulum based on LQR mapping*", 2006 IEEE International Symposium on Intelligent Control.
- [5] Minh Park, Yeoun-Jae Kim, Yeoun-Jae Kim, Ju-Jang Lee, "*Swing-up and LQR stabilization of a rotary inverted pendulum*", Artificial Life and Robotics, 2013.
- [6] Iraj Hassanzadeh, Saleh Mobayen, "*PSO-Based Controller Design for Rotary Inverted Pendulum System*", Journal of Applied Sciences, 2008.

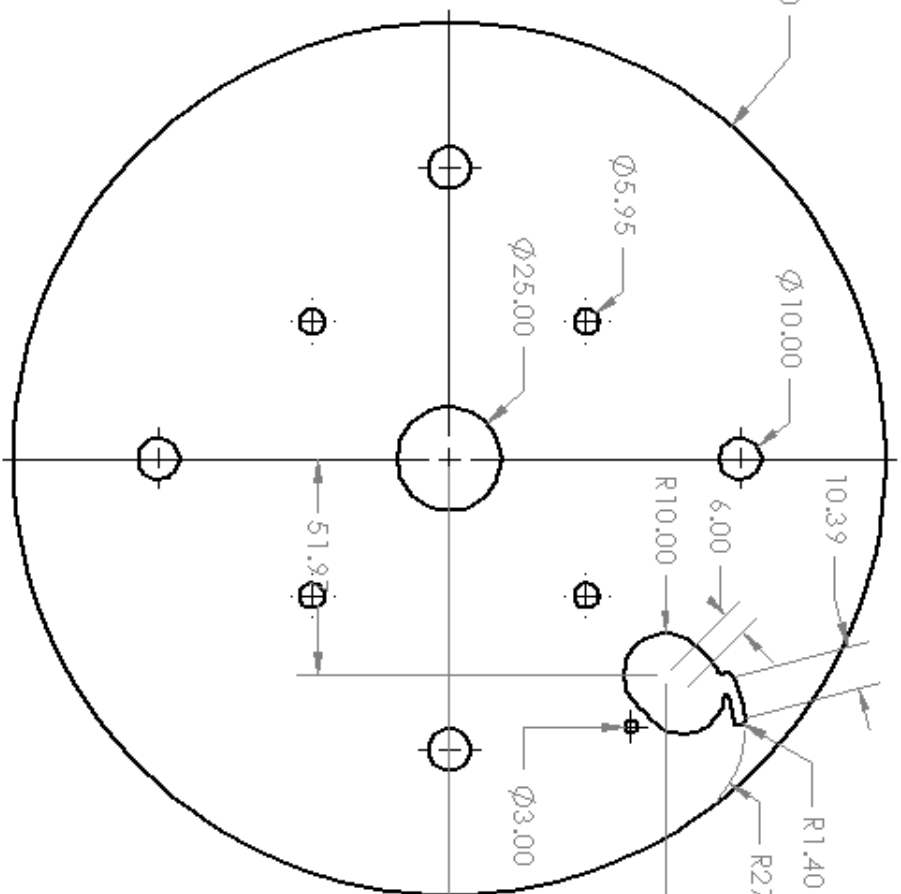
[۷] آ. گشتاسبی، "ساخت و کنترل PID پاندول معکوس دورانی"، ۱۳۹۹.

[۸] م. گرجی، "کنترل پاندول معکوس دورانی به کمک روش‌های هوشمند"، ۱۳۹۹.

پیوست‌ها

پیوست ۱: نقشه کارگاهی قطعات تغییر یافته در آونگ معکوس دورانی

—



➤



Top Plate

Weight

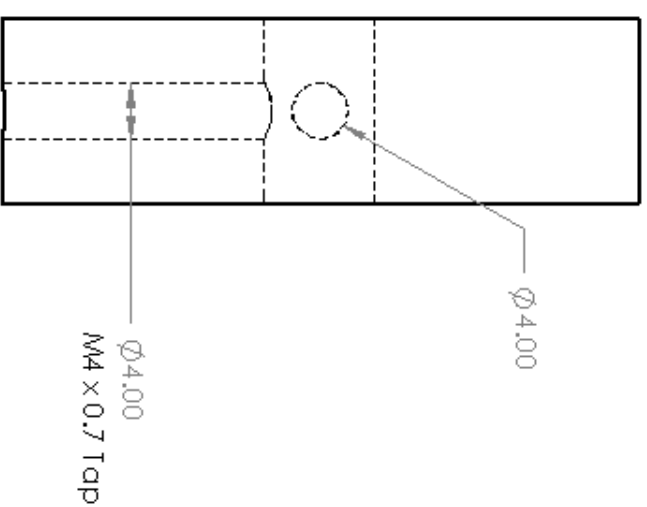
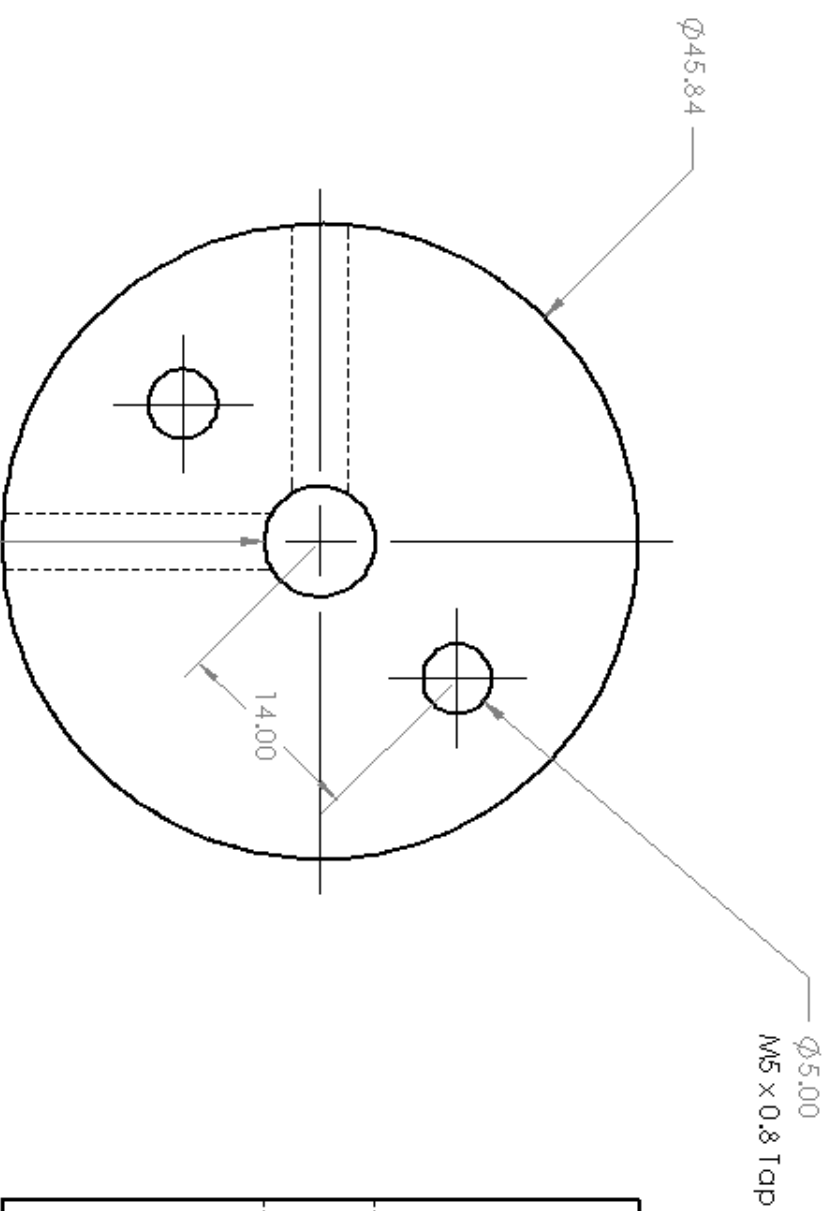
332.41g

SHEET 1 OF 6

—

2

1



A

B

A

B

TITLE:

Motor Pulley
48 teeth 3M

Material

Aluminium

Weight

55.18g

SCALE: 2:1

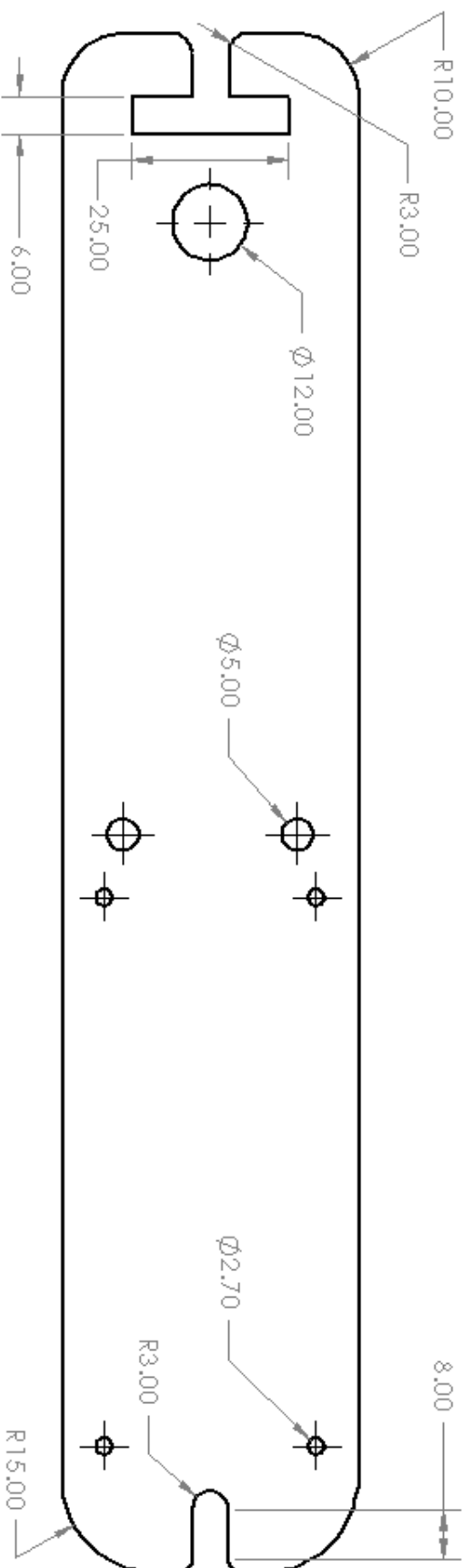
SHEET 2 OF 6

2

1

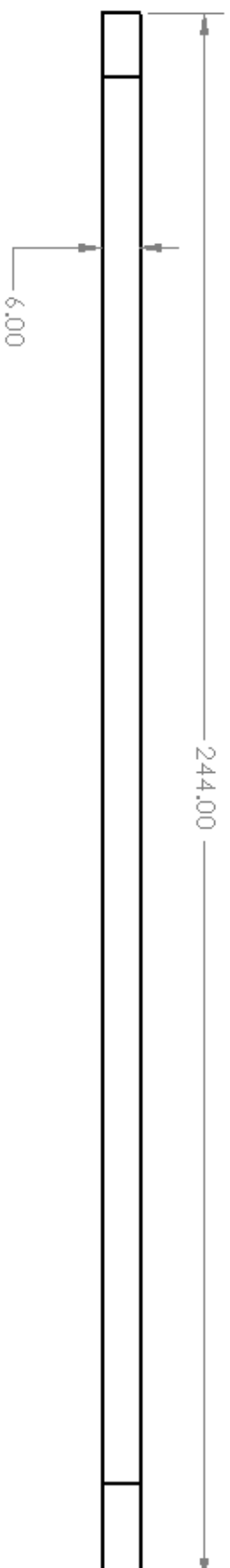
2

1



B

B



A

A

TITLE:
Arm

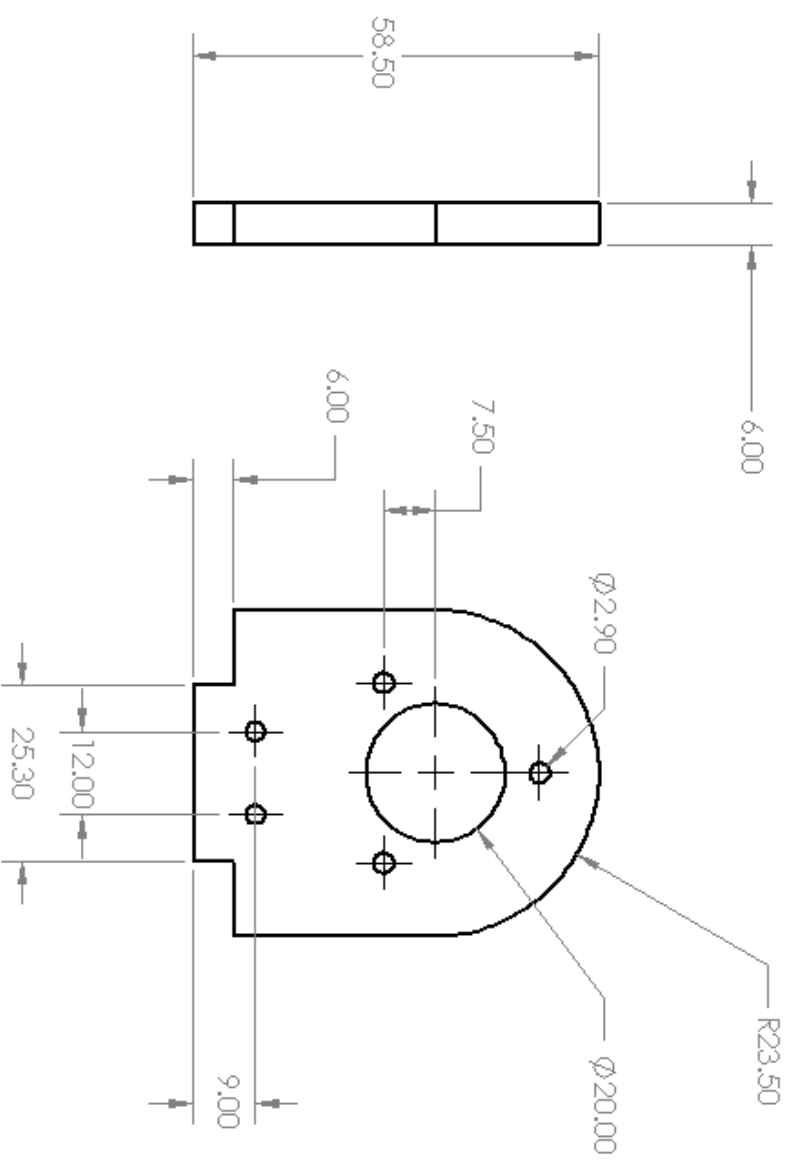
Material	Weight
Plexiglass	76.94g
SCALE: 1:1	SHEET 3 OF 6

2

1

2

1



B

A

B

A

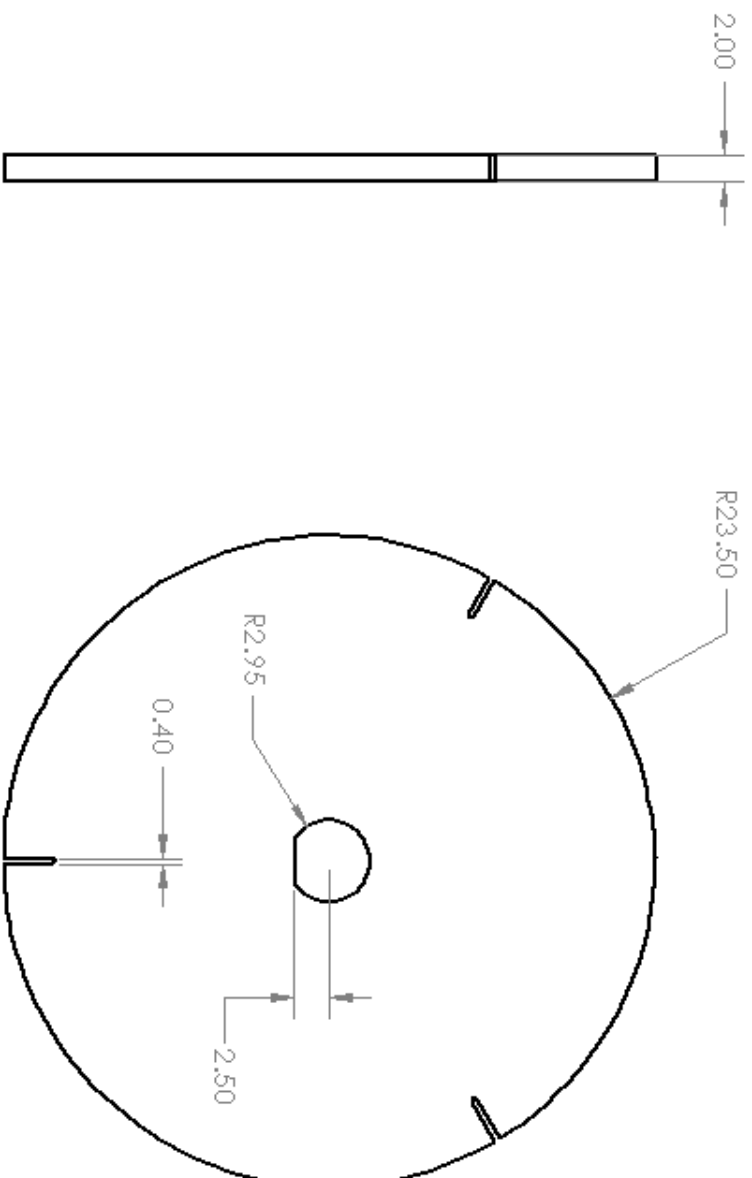
2

1

TITLE: Pendulum Encoder Mounting		
Material	Weight	
Plexiglass	14.42g	
SCALE: 1:1	SHEET 4 OF 6	

2

1



A

B

A

B

2

1

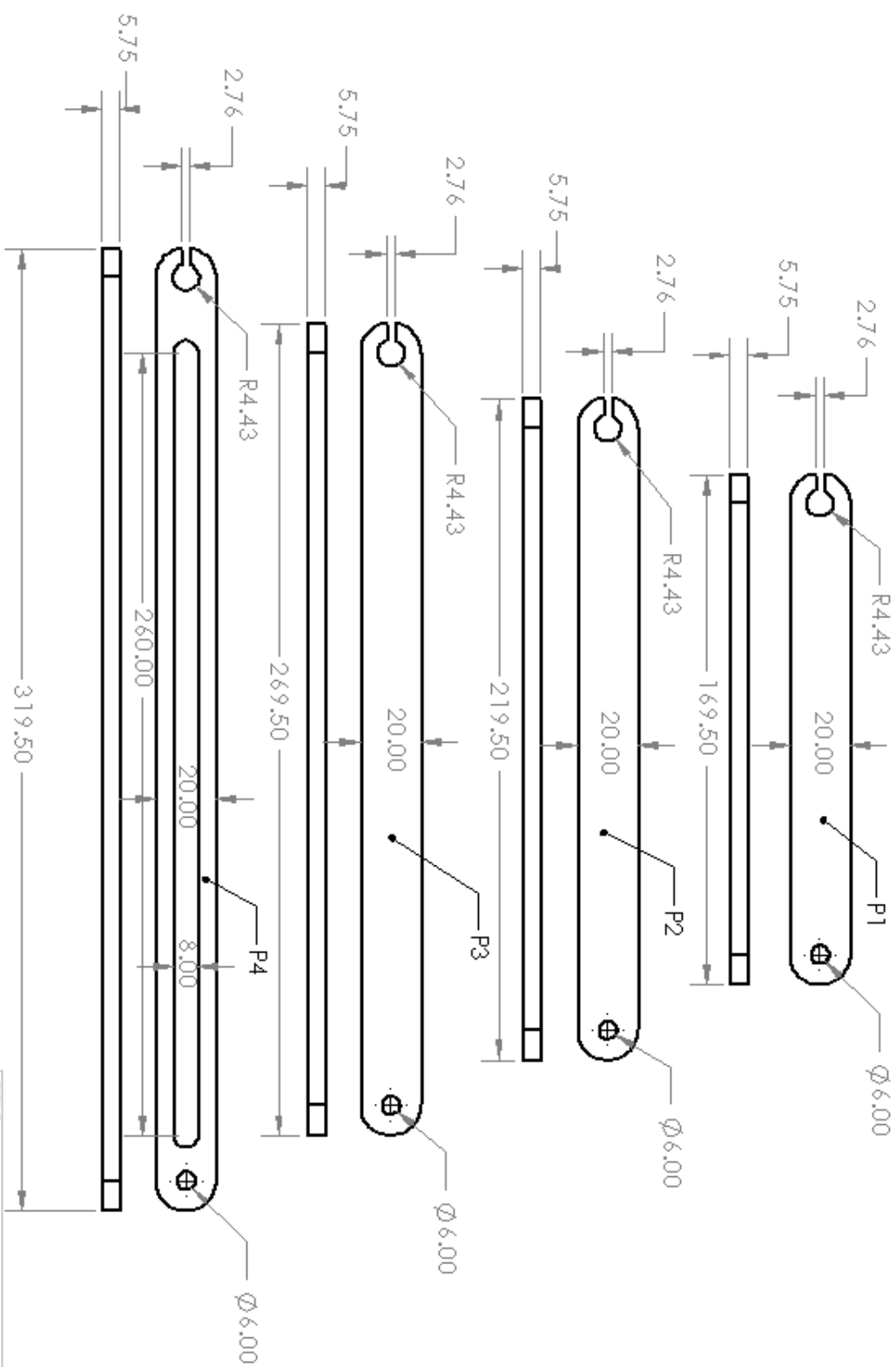
TITLE: Encoder Calibration Disk		
Material	Weight	
Plexiglass	4.02g	
SCALE: 2:1	SHEET 5 OF 6	

2

1

A

B



A

B

Pendulums

Material

Plexiglass

Weights

P1: 21.77g

P2: 28.55g

P3: 35.34g

P4: 27.67g

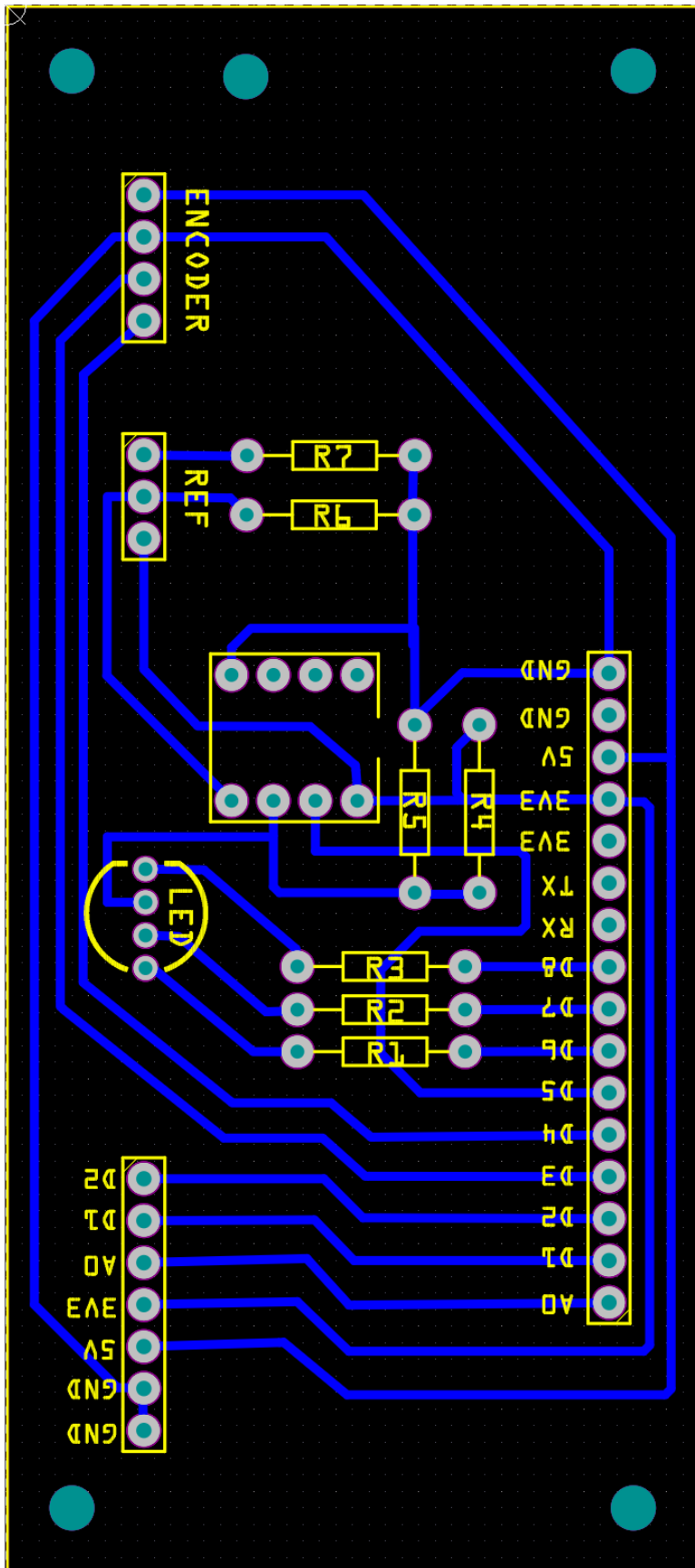
SCALE: 1:2

SHEET 6 OF 6

2

1

پیوست ۲: تصاویر PCBهای طراحی شده در محیط Altium Designer





University of Tehran
College of Engineering
School of Mechanical
Engineering



Control of the Rotary Inverted Pendulum Using Reinforcement Learning Algorithms

**Senior Design Project
Mechanical Engineering B.Sc.**

**Autor: Arash Hatefi
Supervisor: Dr. Masoud Shariatpanahi**

February 2021