



دانشگاه صنعتی امیرکبیر
(پلی تکنیک تهران)

دانشکده مهندسی برق

کنترل و عیب‌یابی هوشمند

پیاده‌سازی روش‌های هوشمند برای کنترل و عیب‌یابی سیستم تعلیق خودروی پیشرفته

نگارش

علی نظامی وند چگینی

استاد

دکتر فرزانه عبدالمهی

بهار 1404

چکیده

مقدمه

سیستم‌های تعلیق نقش مهمی در پایداری، راحتی سواری و عملکرد هندلینگ خودرو دارند. سیستم تعلیق چهارگانه (Quad Car Suspension) یکی از بخش‌های مهم در مهندسی خودرو است که بررسی روش‌های مختلف آن می‌تواند به بهبود دینامیک کلی وسیله نقلیه کمک کند.

در این گزارشکار، مدل‌های مختلفی برای شناسایی و کنترل این سیستم مورد بررسی و آزمایش قرار می‌گیرند. روش‌های متنوعی از جمله مدل‌سازی ریاضی، شبیه‌سازی و مقایسه عملکرد سیستم‌های تعلیق مختلف ارزیابی می‌شوند تا تأثیر هر روش بر کیفیت سواری، پایداری و بازدهی سیستم تحلیل گردد. هدف این گزارشکار، ارائه یک تحلیل جامع از روش‌های مختلف سیستم تعلیق و بررسی تأثیر آن‌ها در شرایط مختلف جاده‌ای و بارگذاری است.

واژه‌های کلیدی:

کنترل، عیب‌یابی، سیستم‌های هوشمند، شبکه‌های عصبی، سیستم تعلیق هوشمند خودرو

فصل اول: مقدمه و مدل سازی سیستم.....	1
معرفی سیستم.....	2
مدل سازی سیستم.....	3
معادلات سیستم.....	4
1-1-1- پارامترهای سیستم.....	5
فصل دوم: کنترل مبتنی بر شبکه MLP.....	7
معرفی روش MLP.....	8
طراحی شبکه MLP.....	12
ساختار شبکه عصبی پرسپترون چندلایه (MLP) در این پیاده سازی.....	14
نتایج شبیه سازی.....	17
فصل سوم: کنترل مبتنی بر شبکه RBF.....	20
معرفی شبکه RBF.....	21
ساختار کنترلر RBF.....	21
ساختار شناسایی کننده RBF.....	24
سیستم حلقه باز.....	26
سیستم حلقه بسته.....	30
نتایج شبیه سازی.....	31
انجام آزمایش به ازای مقادیر مختلف sampling rate.....	33
فصل چهارم: کنترل مبتنی بر شبکه LSTM.....	34
معرفی شبکه LSTM.....	35
ساختار شبکه LSTM.....	35
طراحی مدل و پیاده سازی مدل غیر خطی.....	38
تولید مجموعه داده برای آموزش مدل در سیستم تعلیق فعال خودرو.....	43
مدل پیشنهادی برای شناسایی و رویت سیستم.....	46
نتیجه گیری کلی.....	48
فصل پنجم: کنترل مبتنی بر RL.....	51
معرفی روش کنترل مبتنی بر RL.....	52
معرفی الگوریتم SAC.....	53
توضیح بخش های مختلف برنامه.....	54

65	آموزش شبکه.....
68	نتایج آموزش.....
77	نتیجه گیری کلی.....
79	فصل ششم: جمع بندی و نتیجه گیری.....
80	نتیجه گیری روش های مختلف.....
80	MLP (Perceptron چندلایه).....
80	RBF (شبکه عصبی شعاعی).....
80	LSTM (شبکه های عصبی بازگشتی طولانی مدت کوتاه مدت).....
81	SAC (الگوریتم یادگیری تقویتی Soft Actor-Critic).....
82	مقایسه روش ها.....
82	نتیجه گیری کلی.....
83	منابع و مراجع.....

شکل 1-1 تصویر کلی از سیستم تعلیق خودرو با 4 درجه آزادی	2
شکل 2-1 تصویر کلی از مدل سازی سیستم تعلیق خودرو با 4 درجه آزادی	3
شکل 3-1 مدل سازی سیستم تعلیق پیشرفته	4
شکل 1-2 ساختار کلی شبکه	12
شکل 2-2 خروجی شبکه	17
شکل 3-2 خروجی شبکه به ازای مقادیر مختلف	18
شکل 1-3 ساختار کلی کنترل کننده RBF	22
شکل 2-3 ساختار شناسایی کننده RBF	25
شکل 3-3 خروجی سیستم حلقه باز	27
شکل 4-3 پاسخ حالات سیستم به ورودی پله	28
شکل 5-3 خروجی سیستم نسبت به ورودی سینوسی	30
شکل 6-3 نتایج شبیه سازی حلقه بسته	31
شکل 7-3 تغییرات وزن شناسایی کننده و کنترل کننده	32
شکل 8-3 نتایج شبیه سازی با نرخ نمونه برداری 0.01 ثانیه	33
شکل 1-4 ساختار شبکه عصبی LSTM	37
شکل 2-4 مدلسازی سیستم غیر خطی	38
شکل 3-4 نتایج مدلسازی سیستم (حالات سیستم)	39
شکل 4-4 نتایج مدلسازی سیستم (سیگنال کنترلی)	41
شکل 5-4 خروجی دیتا سازی	44
شکل 6-4 ورودی کنترلی	45
شکل 7-4 نتایج آموزش مدل	47

شکل 5-1	بررسی پاداش و حرکت در طی آموزش	68
شکل 5-2	میزان هزینه شبکه	71
شکل 5-3	critic loss میانگین	72
شکل 5-4	میزان احتمالات اکتور	73
شکل 5-5	alpha temperature تغییر	73
شکل 5-6	critic loss تغییرات	75
شکل 5-7	مانیتور فرکانس هر حالت	76

صفحه

فهرست جداول

جدول 1-1- پارامترهای مورد نیاز برای مدلسازی کوادروتور	5
جدول 2-2- مقایه MLP با روش های کلاسیک	11
جدول 6-1 مقایسه روش های مختلف	82

فصل اول:

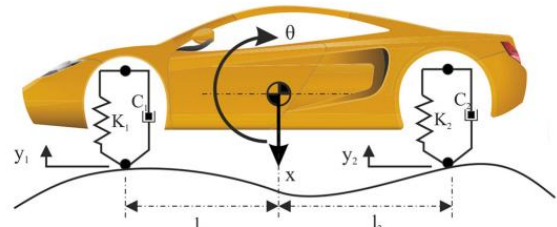
مقدمه و مدل سازی سیستم

معرفی سیستم

سیستم تعلیق خودرو یکی از مهم ترین بخش های آن محسوب می شود که وظیفه اصلی آن، جذب ضربات و نوسانات ناشی از ناهمواری های جاده و حفظ پایداری و کنترل خودرو در شرایط مختلف رانندگی است. سیستم تعلیق چهارگانه (Quad Car Suspension) یک ساختار پیشرفته است که از چهار مجموعه فنر و دمپر مستقل یا نیمه مستقل تشکیل شده و به منظور بهبود کیفیت سواری، کنترل دینامیکی و افزایش ایمنی خودرو طراحی شده است.

این سیستم علاوه بر کاهش لرزش ها و ارتعاشات، نقش مهمی در توزیع بار بین چرخ ها، بهینه سازی چسبندگی تایرها به سطح جاده و بهبود پاسخگویی خودرو به تغییرات مسیر دارد. طراحی و عملکرد آن می تواند بسته به نوع خودرو، کاربری آن و شرایط محیطی متفاوت باشد.

یکی از مهم ترین ویژگی های سیستم تعلیق، رفتار غیر خطی آن است. این غیر خطی بودن ناشی از عواملی مانند سختی متغیر فنرها، رفتار وابسته به سرعت و جابجایی در دمپرها، اصطکاک بین قطعات متحرک، تغییرات زاویه ای و وابستگی نیروهای تایلر به شرایط جاده است. این پیچیدگی ها باعث می شوند که مدل سازی و کنترل این سیستم چالش برانگیز باشد و نیاز به استفاده از روش های پیشرفته شناسایی و کنترل داشته باشد.



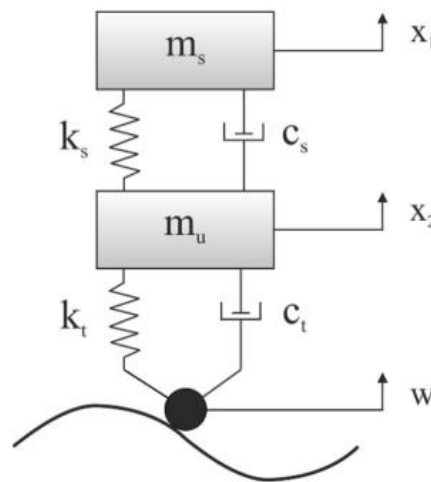
شکل 1-1 تصویر کلی از سیستم تعلیق خودرو با 4 درجه آزادی

مدل سازی سیستم

سیستم تعلیق چهارگانه خودرو یک مدل دینامیکی با چهار درجه آزادی دارد که شامل جرم فنربندی شده، جرم غیرفنربندی شده، فنر و دمپرهای سیستم تعلیق و تأثیر تایرها است. در این مدل، سختی تایر به صورت غیرخطی درجه دوم و سختی فنرهای تعلیق به صورت غیرخطی درجه سوم در نظر گرفته شده است.

معادلات دینامیکی سیستم بر اساس اصل دالامبر استخراج شده اند و نیروهای اعمال شده بر جرمهای سیستم، شامل سختی و میرایی تعلیق و تایر، در نظر گرفته شده اند. این معادلات برای تحلیل و شبیه سازی رفتار سیستم در شرایط مختلف پیاده سازی شده اند.

با توجه به رفتار غیرخطی سیستم، مدل سازی به گونه ای انجام شده که بتواند تأثیرات ناشی از تغییرات سختی فنر و تایر را به طور دقیق شبیه سازی کند. همچنین، از روشهای بهینه سازی برای تنظیم بهینه پارامترهای سیستم تعلیق و بهبود عملکرد آن در زمینه کاهش ارتعاشات و بهبود راحتی سرنشین استفاده شده است.

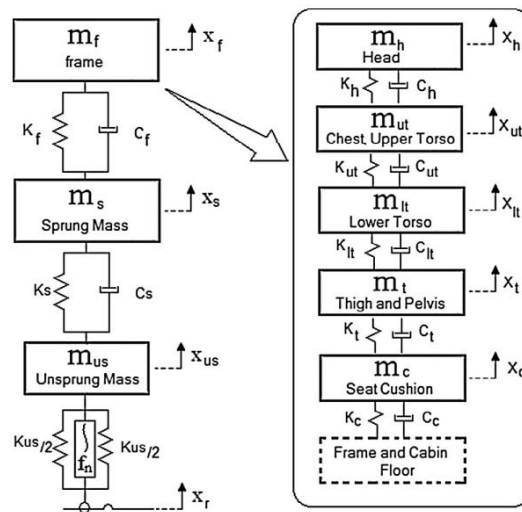


شکل 1-2 تصویر کلی از مدل سازی سیستم تعلیق خودرو با 4 درجه آزادی

معادلات سیستم:

به طور کلی معادلات غیر خطی سیستم به شکل زیر می باشد.

$$\begin{aligned}
 m_{us}\ddot{x}_{us} &= -k_t(x_{us} - x_r) + k_s(x_s - x_{us}) + c_s(\dot{x}_s - \dot{x}_{us}) + k_{tnl}(x_{us} - x_r)^2 + k_{snl}(x_s - x_{us})^3 \\
 m_s\ddot{x}_s &= -k_s(x_s - x_{us}) - c_s(\dot{x}_s - \dot{x}_{us}) - k_{snl}(x_s - x_{us})^3 + k_f(x_f - x_s) + c_f(\dot{x}_f - \dot{x}_s) \\
 m_f\ddot{x}_f &= -k_f(x_f - x_s) - c_f(\dot{x}_f - \dot{x}_s) + k_c(x_c - x_f) + c_c(\dot{x}_c - \dot{x}_f) \\
 m_c\ddot{x}_c &= -k_c(x_c - x_f) - c_c(\dot{x}_c - \dot{x}_f) + k_{tp}(x_t - x_c) + c_{tp}(\dot{x}_t - \dot{x}_c) \\
 m_t\ddot{x}_t &= -k_{tp}(x_t - x_c) - c_{tp}(\dot{x}_t - \dot{x}_c) + k_{lt}(x_{lt} - x_t) + c_{lt}(\dot{x}_{lt} - \dot{x}_t) \\
 m_{lt}\ddot{x}_{lt} &= -k_{lt}(x_{lt} - x_t) - c_{lt}(\dot{x}_{lt} - \dot{x}_t) + k_{ut}(x_{ut} - x_{lt}) + c_{ut}(\dot{x}_{ut} - \dot{x}_{lt}) \\
 m_{ut}\ddot{x}_{ut} &= -k_{ut}(x_{ut} - x_{lt}) - c_{ut}(\dot{x}_{ut} - \dot{x}_{lt}) + k_h(x_h - x_{ut}) + c_h(\dot{x}_h - \dot{x}_{ut}) \\
 m_h\ddot{x}_h &= -k_h(x_h - x_{ut}) - c_h(\dot{x}_h - \dot{x}_{ut})
 \end{aligned}$$



شکل 1-3 مدل سازی سیستم تعلیق پیشرفته

پارامترهای سیستم

پارامترهای مورد نیاز جهت مدل سازی این سیستم تعلیق به صورت زیر می باشد.
جدول 1-1- پارامترهای مورد نیاز برای مدلسازی کوادروتور

Symbol	Description	Unit
A	System matrix	-
AR_h	Amplitude ratio of head RMS acceleration to seat RMS acceleration	-
AR_{ut}	Amplitude ratio of upper torso RMS acceleration to seat RMS acceleration	-
A_{wh}	Frequency-weighted RMS head acceleration	m/s^2
A_{w_spr}	Frequency-weighted RMS sprung mass acceleration	m/s^2
a_{wh}	Frequency-weighted head acceleration	m/s^2
c	Damping coefficient	$N \cdot s/m$
c_{lt}	Lumbar spine damping	$N \cdot s/m$
c_{ut}	Thoracic spine damping	$N \cdot s/m$
c_h	Cervical spine damping	$N \cdot s/m$
f_{obj}	Objective function	-
k	Stiffness	N/m
k_{lt}	Lumbar spine stiffness	N/m
k_{ut}	Thoracic spine stiffness	N/m
k_h	Cervical spine stiffness	N/m
k_{snl}	Nonlinear spring stiffness	N/m^3
k_t	Tire stiffness	N/m
k_{tnl}	Nonlinear tire stiffness	N/m^2
m	Mass	kg
VDV_h	Vibration dose value at the head	$m/s^{1.75}$
x_r	Road profile	m
x, \dot{x}, \ddot{x}	Displacement, velocity, and acceleration	$m, m/s, m/s^2$

• اندیس ها

Subscript	Meaning
s	Sprung mass
us	Unsprung mass
f	Frame
c	Seat cushion
tp	Thigh and pelvis
lt	Lower torso (lumbar spine)
ut	Upper torso
h	Head

➤ مقدار دهی مقادیر:

Parameter	Value	Parameter	Value	Parameter	Value
m_h	5.31	c_h	400	k_h	310,000
m_{ut}	28.49	c_{ut}	4750	k_h	183,000
m_{lt}	8.62	c_{lt}	4585	k_{lt}	162,800
m_t	12.78	c_t	2064	k_t	90,000
m_c	1	c_c	200	k_c	18,000
m_f	15	c_f	830	k_f	31,000
m_s	290	c_s	700	k_s	23,500
m_{us}	40	k_t	190,000	k_{tnl}	$1.5k_t$
		k_{snl}	$100k_s$		

فصل دوم:

کنترل مبتنی بر شبکه MLP

معرفی روش MLP

سیستم تعلیق چهارگانه خودرو یک سیستم پیچیده و غیرخطی است که طراحی کنترل‌کننده‌های کارآمد برای آن نیازمند درک دقیق دینامیک سیستم و واکنش آن نسبت به ورودی‌های مختلف است. از آنجا که مدل‌سازی تحلیلی این سیستم به دلیل وجود نامعینی‌ها، عوامل غیرخطی و تغییرات دینامیکی دشوار است، روش‌های یادگیری ماشین مانند پرسپترون چندلایه (MLP - Multi-Layer Perceptron) گزینه‌ای مناسب برای شناسایی و کنترل این سیستم محسوب می‌شوند.

1) شبکه عصبی MLP و ویژگی‌های آن

MLP یکی از پرکاربردترین شبکه‌های عصبی مصنوعی (ANN) است که شامل چندین لایه متصل به هم شامل لایه ورودی، یک یا چند لایه پنهان و لایه خروجی می‌باشد. ویژگی کلیدی این شبکه، توانایی یادگیری نگاشت‌های غیرخطی پیچیده بین متغیرهای ورودی و خروجی است. هر نورون در لایه‌های پنهان دارای یک تابع فعال‌سازی (مانند سیگموئید، تانژانت هیپربولیک یا ReLU) است که باعث افزایش قدرت یادگیری شبکه در مسائل پیچیده می‌شود.

مزایای استفاده از MLP برای سیستم تعلیق چهارگانه خودرو

- ✓ قابلیت شناسایی رفتارهای غیرخطی سیستم
- ✓ توانایی تعمیم به ورودی‌های جدید و مقابله با نامعینی‌ها
- ✓ امکان استفاده در طراحی کنترل‌کننده‌های تطبیقی و مقاوم
- ✓ کاهش نیاز به مدل‌سازی دقیق تحلیلی

2) شناسایی سیستم تعلیق خودرو با استفاده از MLP

مرحله 1: جمع آوری داده‌ها

در این مرحله، مجموعه‌ای از داده‌های تجربی یا شبیه‌سازی شده از رفتار سیستم تحت ورودی‌های مختلف جمع‌آوری می‌شود. این داده‌ها شامل:

- ورودی‌ها: نیروی جاده، دست‌اندازها، سرعت خودرو و نیروهای تعلیق
- خروجی‌ها: جابه‌جایی چرخ، شتاب سرنشین، نیروی واکنش تعلیق

مرحله 2: آموزش شبکه MLP

- ساختار شبکه: تعداد نرون‌ها و لایه‌های پنهان با استفاده از روش‌های بهینه‌سازی مانند آزمایش و خطا (Trial and Error) یا جستجوی شبکه‌ای (Grid Search) انتخاب می‌شوند.
- الگوریتم آموزش: از الگوریتم‌های مبتنی بر پس‌انتشار خطا (Backpropagation) مانند بهینه‌سازی گرادیان نزولی (SGD) یا آدام (Adam) برای تنظیم وزن‌های شبکه استفاده می‌شود.
- داده‌های آموزشی و تست: داده‌ها به دو بخش آموزشی و تست تقسیم شده و مدل روی داده‌های آموزشی یاد می‌گیرد، سپس روی داده‌های تست ارزیابی می‌شود.

مرحله 3: ارزیابی مدل شناسایی

پس از آموزش، مدل باید بتواند رفتار سیستم تعلیق را به‌طور دقیق شبیه‌سازی کند. برای ارزیابی عملکرد آن از معیارهایی مانند:

- میانگین مربع خطا (MSE)
- ضریب همبستگی (R^2)

استفاده می‌شود. در صورت عدم دقت کافی، ساختار شبکه یا پارامترهای یادگیری تنظیم می‌شوند.

3) طراحی کنترل کننده مبتنی بر MLP برای سیستم تعلیق خودرو

کنترل کننده های مبتنی بر MLP می توانند به دو روش پیاده سازی شوند:

الف) کنترل مستقیم (Direct Control)

در این روش، شبکه عصبی مستقیماً یک نگاشت بین ورودی های سیستم (مانند دست انداز جاده و نیروهای تعلیق) و سیگنال های کنترلی (مانند گشتاور یا نیروی اعمالی به کمک فنرها) یاد می گیرد. این روش نیازمند حجم زیادی از داده های آموزشی برای یادگیری یک سیاست کنترلی کارآمد است.

ب) کنترل غیرمستقیم (Indirect Control)

در این رویکرد، ابتدا یک مدل دینامیکی از سیستم با استفاده از شبکه MLP به دست می آید، سپس از این مدل برای طراحی یک کنترل کننده مناسب (مانند کنترل فیدبک تطبیقی یا کنترل پیش بین مدل ((MPC)) استفاده می شود.

الگوریتم های کنترل رایج در کنار MLP

✓ **کنترل پیش بین مدل (MPC):** از مدل آموزش دیده شده برای پیش بینی رفتار آینده سیستم و بهینه سازی سیگنال های کنترلی استفاده می شود.

✓ **کنترل تطبیقی:** از شبکه MLP برای به روز رسانی پارامترهای کنترلی در شرایط نامشخص استفاده می شود.

✓ **کنترل مقاوم (Robust Control):** ترکیب MLP با روش های مقاوم برای کاهش حساسیت به تغییرات پارامترهای سیستم.

4) مقایسه MLP با روش‌های کلاسیک کنترل

جدول 2-2- مقایسه MLP با روش‌های کلاسیک

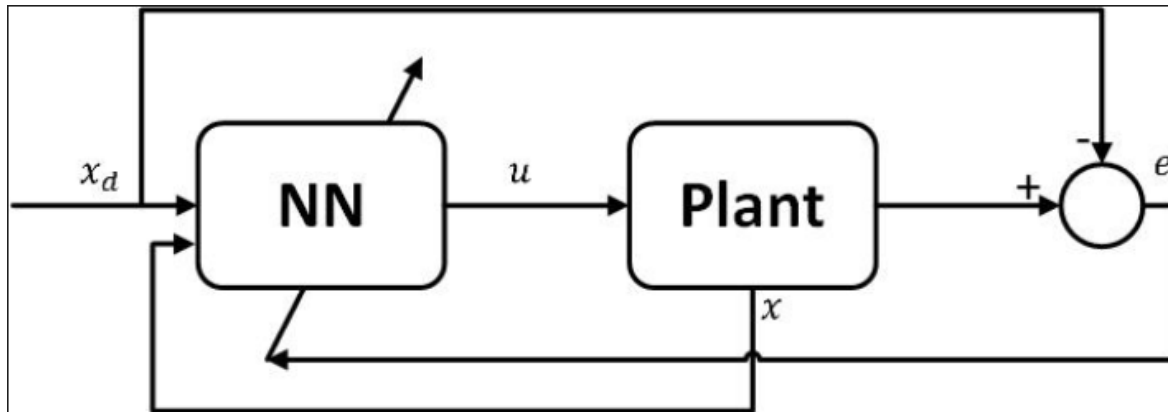
ویژگی	کنترل کلاسیک (PID, LQR)	کنترل مبتنی بر MLP
نیاز به مدل ریاضی	بله، نیاز به مدل دقیق دارد	خیر، می‌تواند بدون مدل تحلیلی یاد بگیرد
قابلیت انطباق	ضعیف در برابر تغییرات سیستم	قوی، می‌تواند به صورت آنلاین به روزرسانی شود
شناسایی سیستم	مبتنی بر معادلات دینامیکی	مبتنی بر داده‌های تجربی
پیچیدگی محاسباتی	کم، مناسب برای پیاده‌سازی سریع	بالا، نیازمند توان پردازشی بیشتر
عملکرد در سیستم‌های غیرخطی	ضعیف یا نیازمند خطی‌سازی	قوی، بدون نیاز به خطی‌سازی

5) چالش‌های استفاده از MLP در کنترل و شناسایی سیستم تعلیق

- ✓ نیاز به داده‌های آموزشی گسترده برای دستیابی به دقت بالا
- ✓ پیچیدگی در انتخاب معماری بهینه شبکه و تنظیم پارامترهای یادگیری
- ✓ هزینه محاسباتی بالاتر نسبت به روش‌های کلاسیک

طراحی شبکه MLP

ساختار کلی که ما در این پروژه داریم به شکل زیر می باشد.



شکل 1-2 ساختار کلی شبکه

در این پیاده سازی، یک شبکه عصبی پرسپترون چندلایه (MLP) طراحی شده است که می تواند در مسائل شناسایی و کنترل سیستم های مختلف مورد استفاده قرار گیرد. این نوع شبکه عصبی از مجموعه ای از نورون ها در چندین لایه تشکیل شده است که با استفاده از ارتباطات وزنی بین آن ها، فرآیند یادگیری را انجام می دهد. ساختار شبکه شامل لایه ورودی، چندین لایه پنهان و یک لایه خروجی است که هر کدام نقش مهمی در پردازش داده ها و یادگیری الگوهای مختلف ایفا می کنند.

این مدل بر اساس اصول یادگیری ماشین و روش های مبتنی بر شبکه های عصبی طراحی شده است. در این روش، ورودی های شبکه پس از عبور از لایه های مختلف، با استفاده از توابع فعال سازی پردازش شده و در نهایت خروجی مناسب تولید می شود. توابع فعال سازی نقش مهمی در تعیین پاسخ شبکه به ورودی های مختلف دارند و به آن امکان می دهند تا روابط پیچیده بین داده ها را بیاموزد. انتخاب تابع فعال سازی مناسب تأثیر بسزایی در عملکرد شبکه دارد و می تواند رفتار غیرخطی سیستم را به خوبی مدل کند.

در فرآیند یادگیری، ابتدا وزن‌ها و بایاس‌های شبکه به صورت مقداردهی اولیه تنظیم می‌شوند. سپس از طریق فرآیند انتشار رو به جلو، مقدار خروجی شبکه محاسبه شده و با مقدار هدف مقایسه می‌شود. اختلاف بین خروجی واقعی و مقدار مطلوب، خطا را تشکیل می‌دهد که باید در فرآیند یادگیری کاهش یابد. برای این منظور، روش پس انتشار خطا استفاده می‌شود که از طریق محاسبه گرادیان‌ها، وزن‌های شبکه را به روزرسانی کرده و باعث بهبود عملکرد شبکه در تطبیق با داده‌های ورودی می‌شود.

این روش یادگیری، مبتنی بر بهینه‌سازی وزن‌ها از طریق الگوریتم‌های گرادیان نزولی انجام می‌شود. نرخ یادگیری در این فرآیند نقش مهمی ایفا می‌کند و تعیین مقدار مناسب آن می‌تواند تأثیر قابل توجهی بر سرعت همگرایی و دقت نهایی مدل داشته باشد. در این پیاده‌سازی، پس از محاسبه گرادیان‌ها، وزن‌ها و بایاس‌های شبکه با اعمال تغییرات کوچک به سمت مقادیر بهینه هدایت می‌شوند. این فرآیند به صورت تکراری انجام شده و تا زمانی که شبکه به دقت مطلوبی برسد ادامه می‌یابد.

شبکه‌های عصبی پرسپترون چندلایه به دلیل توانایی در یادگیری الگوهای پیچیده و مدل‌سازی روابط غیرخطی، کاربردهای گسترده‌ای در زمینه‌های مختلف دارند. از جمله کاربردهای آن می‌توان به شناسایی سیستم‌های غیرخطی، کنترل تطبیقی، پردازش سیگنال، تشخیص الگو، و پیش‌بینی داده‌ها اشاره کرد. این مدل می‌تواند در سیستم‌هایی که دارای رفتارهای پیچیده و غیرخطی هستند، به عنوان یک ابزار کارآمد برای تحلیل و کنترل مورد استفاده قرار گیرد.

در این پیاده‌سازی، ساختار شبکه به گونه‌ای طراحی شده است که انعطاف‌پذیری بالایی در تنظیم تعداد لایه‌ها و نوروں‌های هر لایه دارد. این ویژگی امکان تنظیم معماری شبکه بر اساس نیازهای خاص هر مسئله را فراهم می‌کند. همچنین از روش‌های مختلف فعال‌سازی برای بهبود عملکرد شبکه در شرایط مختلف استفاده شده است.

در نظر داریم که روابط استفاده شده در این پیاده‌سازی به شکل زیر می‌باشد:

$$\dot{x} = A_c x + f(x) + g(x)u, \quad \text{where } f(x) = f_0(x) - A_c x$$

$$e = x - x_d, \quad J = \frac{1}{2}(e^T e) + u^T R u$$

$$\frac{\partial J}{\partial e} = e$$

$$\text{considering static approximation: } \dot{x} = 0 \rightarrow \frac{\partial e}{\partial u} = -A_c^{-1} g(x) \approx -A_c^{-1} I_{n \times 1}$$

$$\frac{\partial J}{\partial u} = \frac{\partial J}{\partial e} \frac{\partial e}{\partial u} + \frac{\partial (u^T R u)}{\partial u} = -e A_c^{-1} I_{n \times 1} + u^T R$$

ساختار شبکه عصبی پرسپترون چندلایه (MLP) در این پیاده‌سازی

شبکه عصبی پرسپترون چندلایه (MLP) طراحی شده در این کد، برای کنترل یک سیستم تعلیق طراحی شده است. این شبکه دارای یک ساختار استاندارد با لایه ورودی، دو لایه پنهان و یک لایه خروجی است. هدف از این طراحی، استفاده از یک کنترل کننده عصبی برای کاهش نوسانات سیستم تعلیق و بهبود پایداری آن است.

۱. معماری کلی شبکه

این شبکه از چندین لایه تشکیل شده است که هر یک دارای نقش مشخصی در فرآیند یادگیری و کنترل دارند:

- **لایه ورودی:** شامل ۴ نورون است که نمایانگر وضعیت سیستم تعلیق می‌باشد. این ورودی‌ها شامل:

- جابه‌جایی جرم اصلی (x_1)

- جابه‌جایی چرخ (x_2)

○ سرعت جرم اصلی (v_1)

○ سرعت چرخ (v_2)

این مقادیر وضعیت لحظه‌ای سیستم را نشان می‌دهند و ورودی شبکه عصبی را تشکیل می‌دهند.

- لایه‌های پنهان: دو لایه پنهان با ۳۲ نورون در هر لایه در نظر گرفته شده‌اند. این لایه‌ها وظیفه استخراج ویژگی‌های پیچیده از ورودی‌ها را بر عهده دارند. در این لایه‌ها از تابع فعال‌سازی ReLU استفاده شده است که یک انتخاب مناسب برای مدل‌سازی روابط غیرخطی در سیستم‌های دینامیکی است.

- لایه خروجی: این لایه دارای ۱ نورون است که مقدار سیگنال کنترلی (u) را تولید می‌کند. این مقدار، نیرویی است که توسط سیستم کنترل به سیستم تعلیق اعمال می‌شود تا نوسانات را کاهش دهد. برای این لایه، تابع فعال‌سازی خطی (Linear) استفاده شده است که مقدار کنترلی را بدون محدودیت مستقیماً ارائه می‌دهد.

۲. مکانیسم یادگیری و بهینه‌سازی

- انتشار رو به جلو: (Forward Propagation)

ورودی‌ها از طریق وزن‌ها و بایاس‌های لایه‌های مختلف عبور کرده و با استفاده از توابع فعال‌سازی پردازش می‌شوند تا خروجی نهایی (نیروی کنترلی u) تولید شود.

- محاسبه خطا:

خروجی شبکه با مقدار مطلوب مقایسه شده و خطا محاسبه می‌شود. در اینجا، مقدار مطلوب برابر با صفر در نظر گرفته شده است که نشان‌دهنده هدف کاهش نوسانات سیستم تعلیق است.

- پس انتشار خطا: (Backpropagation)

خطا از طریق الگوریتم پس انتشار به لایه‌های قبلی منتقل شده و وزن‌های شبکه بر اساس آن اصلاح می‌شوند.

• به روزرسانی وزن‌ها:

از طریق روش **گرادیان نزولی**، وزن‌ها و بایاس‌ها با استفاده از نرخ یادگیری مشخص، تنظیم می‌شوند تا شبکه به مرور زمان پاسخ بهتری ارائه دهد. در این پیاده‌سازی، نرخ یادگیری برابر با ۰.۱ در نظر گرفته شده و تعداد مراحل به‌روزرسانی در هر گام یادگیری ۵۰ بار است که به همگرایی سریع‌تر مدل کمک می‌کند.

۳. تعامل شبکه عصبی با سیستم تعلیق

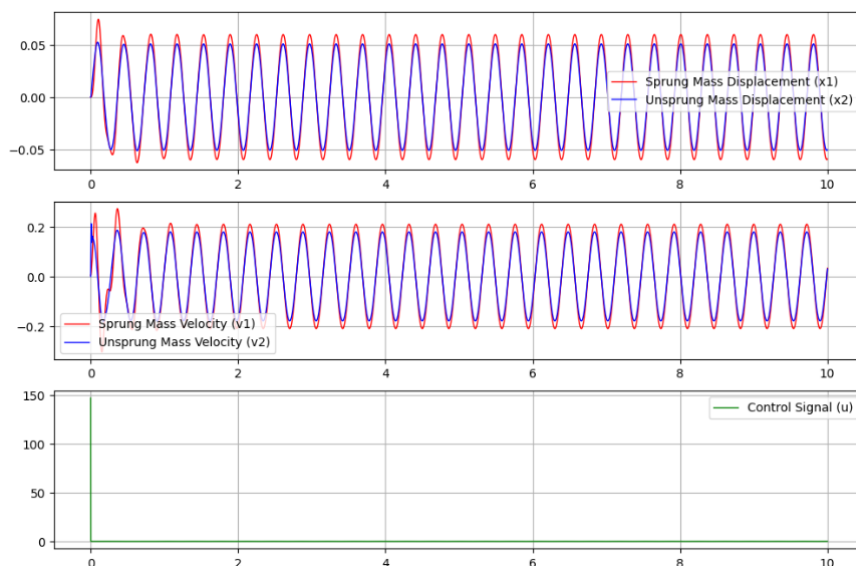
- در هر لحظه، وضعیت سیستم تعلیق اندازه‌گیری شده و به‌عنوان ورودی به شبکه عصبی داده می‌شود.
- شبکه عصبی مقدار نیروی کنترلی (u) را پیش‌بینی کرده و به سیستم تعلیق اعمال می‌کند.
- سپس، سیستم تعلیق با این نیروی جدید، وضعیت خود را به‌روزرسانی کرده و مقدار جدید جابه‌جایی و سرعت را ارائه می‌دهد.
- این فرآیند به‌صورت تکراری اجرا شده و شبکه به‌مرور عملکرد خود را بهینه می‌کند.

۴. کاربردها و مزایای این شبکه

- **توانایی یادگیری روابط پیچیده و غیرخطی:** سیستم تعلیق دارای رفتار غیرخطی است و استفاده از یک کنترل‌کننده عصبی به‌جای روش‌های کلاسیک می‌تواند عملکرد بهتری ارائه دهد.
- **بهبود عملکرد سیستم تعلیق:** این شبکه باعث کاهش نوسانات و بهبود پایداری خودرو در مواجهه با ناهمواری‌های جاده‌ای می‌شود.
- **انعطاف‌پذیری بالا:** معماری شبکه می‌تواند به‌راحتی تغییر کند تا برای سیستم‌های دینامیکی مختلف بهینه شود.

نتایج شبیه سازی

با توجه به اینکه تعداد دور بروز رسانی وزن ها را برابر 50 و زمان نمونه برداری را 0.01 ثانیه و نرخ یادگیری را 0.1 در نظر گرفتیم خروجی شبکه به شکل زیر شده است.



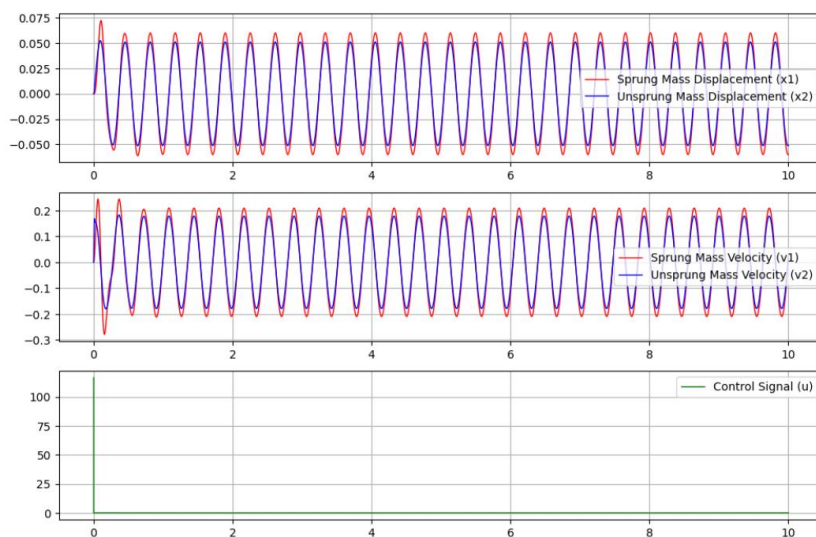
شکل 2-2 خروجی شبکه

همانطور که مشاهده میکنید با خطای بسیار کمی سیگنال مرجع دنبال شده است.

تاثیر هایپر پارامتر ها بر نتیجه tracking :

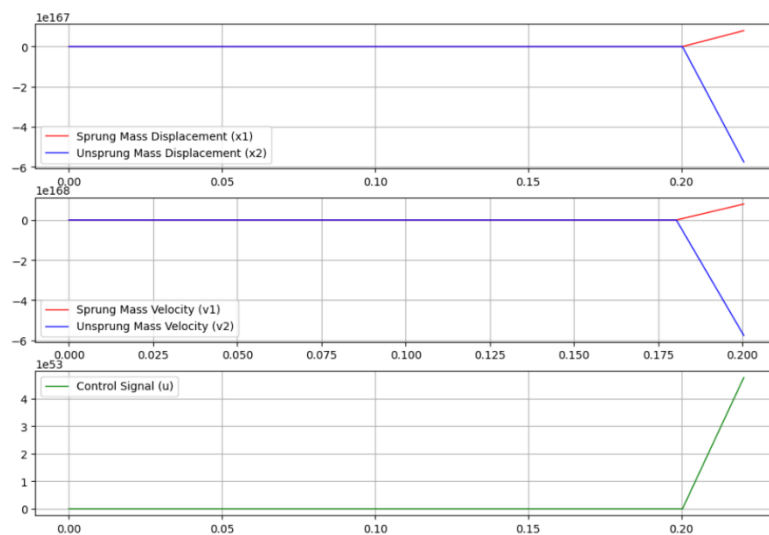
$L_r = 0.01$

$T_d = 0.001$



$L_r = 0.01$

$T_d = 0.1$



شکل 2-3 خروجی شبکه به ازای مقادیر مختلف

همانطور که مشاهده میکنید با افزایش sampling rate دقت افزایش می یابد اما حجم محاسبات شبکه به شدت بالا می رود اما از طرفی با کاهش این مقدار شبکه ممکن است حتی رو به ناپایداری برود .

لذا در این بخش ما مقادیر اولیه را به عنوان بهینه ترین حالت در نظر میگیریم.

فصل سوم:

کنترل مبتنی بر شبکه RBF

معرفی شبکه RBF

در این بخش، از شبکه عصبی پایه شعاعی (RBF) به عنوان یک ابزار قدرتمند برای شناسایی و کنترل سیستم استفاده شده است. این شبکه با بهره‌گیری از توابع پایه شعاعی، قادر است الگوهای پیچیده و غیرخطی را به خوبی مدل‌سازی کند. در فرآیند شناسایی سیستم، شبکه RBF ورودی‌ها و خروجی‌های سیستم را دریافت کرده و با یادگیری روابط میان آن‌ها، مدلی تقریبی از رفتار سیستم ارائه می‌دهد. این ویژگی به ویژه در سیستم‌های دینامیکی غیرخطی که مدل‌سازی ریاضی آن‌ها دشوار است، اهمیت بالایی دارد.

در زمینه کنترل، شبکه RBF به عنوان یک کنترل‌کننده هوشمند به کار گرفته می‌شود که با استفاده از داده‌های به دست آمده از شناسایی سیستم، ورودی‌های کنترلی مناسبی تولید می‌کند تا خروجی سیستم به مقدار مطلوب نزدیک شود. این روش نه تنها باعث افزایش دقت کنترل می‌شود، بلکه توانایی سازگاری با تغییرات دینامیکی سیستم را نیز دارد. به دلیل ساختار ساده و توانایی بالای شبکه‌های RBF در تقریب توابع پیچیده، این روش کنترل و شناسایی می‌تواند عملکرد بهتری نسبت به روش‌های کلاسیک ارائه دهد و در بسیاری از کاربردهای مهندسی مورد استفاده قرار گیرد.

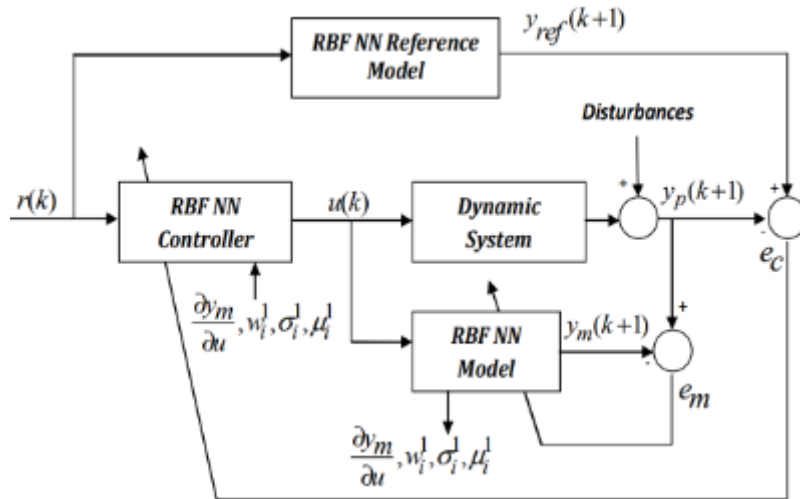
ساختار کنترلر RBF

کنترل‌کننده مبتنی بر شبکه عصبی پایه شعاعی (RBF) یکی از روش‌های قدرتمند در کنترل سیستم‌های دینامیکی است که از قابلیت‌های یادگیری شبکه‌های عصبی برای تخمین و تنظیم ورودی‌های کنترلی استفاده می‌کند. این کنترل‌کننده با استفاده از توابع پایه شعاعی، داده‌های ورودی را به یک فضای ویژگی تبدیل کرده و از طریق یک ساختار ساده اما مؤثر، ورودی کنترلی مناسب را تولید می‌کند.

در طراحی این کنترل‌کننده، لایه پنهان شامل تعدادی نورون با مراکز و پهنای مشخصی است که تعیین‌کننده پاسخ شبکه به ورودی‌های مختلف هستند. در مرحله یادگیری، این پارامترها به گونه‌ای تنظیم می‌شوند که شبکه بتواند الگوی مناسبی از رفتار سیستم را مدل کند. فرآیند یادگیری شامل تنظیم وزن‌های خروجی، مراکز توابع پایه و پارامترهای گوسی توابع پایه است که از طریق قوانین یادگیری مبتنی بر گرادیان و سیگنال خطا بهینه‌سازی می‌شوند.

یکی از ویژگی‌های مهم این کنترل کننده، **نرمال سازی ورودی‌ها** است که باعث بهبود عملکرد شبکه و جلوگیری از ناپایداری در یادگیری می‌شود. سپس، با استفاده از یک مکانیزم انتشار پیشرو، مقدار خروجی شبکه محاسبه می‌شود. این مقدار نشان‌دهنده پاسخ سیستم عصبی به ورودی داده شده است. علاوه بر این، در فرآیند یادگیری و به‌روزرسانی وزن‌ها، از یک تابع خطا برای تنظیم پارامترهای شبکه استفاده می‌شود که منجر به افزایش دقت و بهبود همگرایی در عملکرد کنترل کننده می‌شود.

به‌طور کلی، این کنترل کننده با داشتن ساختار ساده اما انعطاف‌پذیر، توانایی بالایی در مقابله با **سیستم‌های غیرخطی و پیچیده** دارد. استفاده از آن در کاربردهای کنترلی می‌تواند موجب بهبود عملکرد سیستم، کاهش خطا و افزایش تطبیق‌پذیری با تغییرات محیطی شود.



شکل 3-1 ساختار کلی کنترل کننده RBF

روابط کنترلر RBF:

$$\psi_i = \exp \left(-\frac{\|x_{n_c}(k) - c_i\|^2}{2\hat{\sigma}_i^2} \right)$$

$$u(k) = \sum_{i=1}^{n_2} v_i \psi_i$$

$$e_c = y_{ref}(k) - y(k)$$

$$v_i(k+1) = v_i(k) + \eta_{c_1} e_c \psi_i \sum_{i=1}^{n_1} w_i \phi_i \left(-\frac{u(k) - \mu_{i1}}{\sigma_i^2} \right)$$

$$\hat{\sigma}_i(k+1) = \hat{\sigma}_i(k) + \eta_{c_2} e_c v_i \frac{\|x_{n_c}(k) - c_i\|^2}{\hat{\sigma}_i^3} \psi_i \sum_{i=1}^{n_1} w_i \phi_i \left(-\frac{u(k) - \mu_{i1}}{\sigma_i^2} \right)$$

$$c_i(k+1) = c_i(k) + \eta_{c_3} e_c v_i \frac{x_{n_c}(k) - c_i}{\hat{\sigma}_i^2} \psi_i \sum_{i=1}^{n_1} w_i \phi_i \left(-\frac{u(k) - \mu_{i1}}{\sigma_i^2} \right)$$

$$0 \leq \eta_{c_i} \leq 1$$

ساختار شناسایی کننده RBF

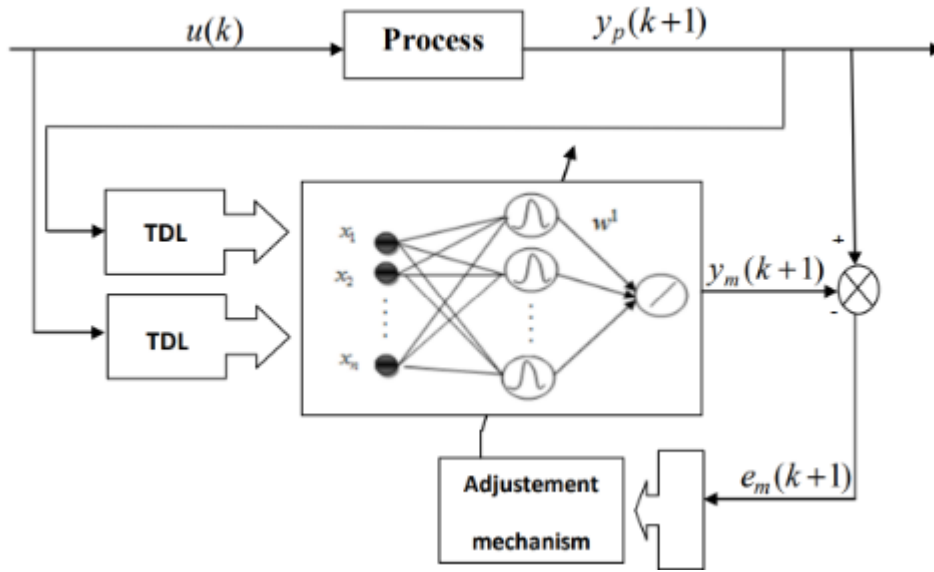
شناسایی کننده مبتنی بر شبکه عصبی پایه شعاعی (RBF) یک مدل هوشمند برای تخمین و شناسایی رفتار سیستم‌های دینامیکی است که بر اساس ویژگی‌های توابع پایه شعاعی عمل می‌کند. این شبکه قادر است روابط غیرخطی پیچیده را بین ورودی‌ها و خروجی‌های یک سیستم مدل کند و یک تخمین دقیق از رفتار سیستم ارائه دهد.

در این شناسایی کننده، از یک مجموعه نورون‌های RBF در لایه پنهان استفاده می‌شود. هر نورون دارای مرکز (μ) و پهنای (σ) مخصوص به خود است که میزان تأثیر آن بر ورودی‌ها را تعیین می‌کند. این مقادیر در طول فرآیند یادگیری تنظیم می‌شوند تا تطابق بهتری با داده‌های ورودی و خروجی سیستم برقرار شود. لایه خروجی شبکه نیز شامل وزن‌های W است که خروجی نهایی را با ترکیب پاسخ نورون‌های پنهان ایجاد می‌کند.

ورودی‌های شناسایی کننده ابتدا نرمال‌سازی می‌شوند تا تأثیر مقیاس متغیرهای مختلف کاهش یابد. سپس، مقدار ϕ ، که نشان‌دهنده فعال‌سازی هر نورون نسبت به ورودی است، محاسبه می‌شود. این مقدار بر اساس فاصله ورودی از مراکز نورون‌ها و تابع گوسی تعیین می‌شود. در نهایت، خروجی شبکه از طریق ترکیب وزن‌دار مقادیر ϕ تولید می‌شود.

فرآیند یادگیری شامل به‌روزرسانی وزن‌ها، مراکز و پهنای نورون‌ها بر اساس سیگنال خطا است. سیگنال خطا، اختلاف بین مقدار واقعی و مقدار تخمین‌زده شده توسط شبکه است. با استفاده از قواعد یادگیری مبتنی بر گرادیان، پارامترهای شبکه تنظیم شده و عملکرد آن در طول زمان بهبود می‌یابد.

این شناسایی کننده با قابلیت تعمیم بالا و توانایی یادگیری رفتارهای غیرخطی، به‌طور گسترده در مسائل کنترل تطبیقی، تخمین پارامترهای سیستم و پردازش سیگنال استفاده می‌شود. ساختار انعطاف‌پذیر و توانایی همگرایی سریع از ویژگی‌های بارز این مدل محسوب می‌شود که آن را به یک ابزار قدرتمند در تحلیل و مدل‌سازی سیستم‌های پیچیده تبدیل کرده است.



شکل 2-3 ساختار شناسایی کننده RBF

روابط شناسایی کننده RBF:

$$\phi_i = \exp \left(-\frac{\|x_{n_{id}}(k) - \mu_i\|^2}{2\sigma_i^2} \right)$$

$$y_m = \sum_{i=1}^{n_1} w_i \phi_i$$

$$e_m = y(k) - y_m$$

$$w_i(k+1) = w_i(k) + \eta_1 e_m \phi_i$$

$$\sigma_i(k+1) = \sigma_i(k) + \eta_2 w_i \frac{\|x_{n_{id}}(k) - \mu_i\|^2}{\sigma_i^3} e_m \phi_i$$

$$\mu_i(k+1) = \mu_i(k) + \eta_3 w_i \frac{x_{n_{id}}(k) - \mu_i}{\sigma_i^2} e_m \phi_i$$

$$0 \leq \eta_i \leq 1$$

سیستم حلقه باز

شبیه‌سازی حلقه باز (**Open-Loop Simulation**) یک روش مهم برای بررسی رفتار سیستم‌های دینامیکی بدون اعمال کنترل‌کننده خارجی است. این روش به ما امکان می‌دهد تا واکنش سیستم را تنها بر اساس ورودی‌های مشخص و معادلات دینامیکی آن تحلیل کنیم.

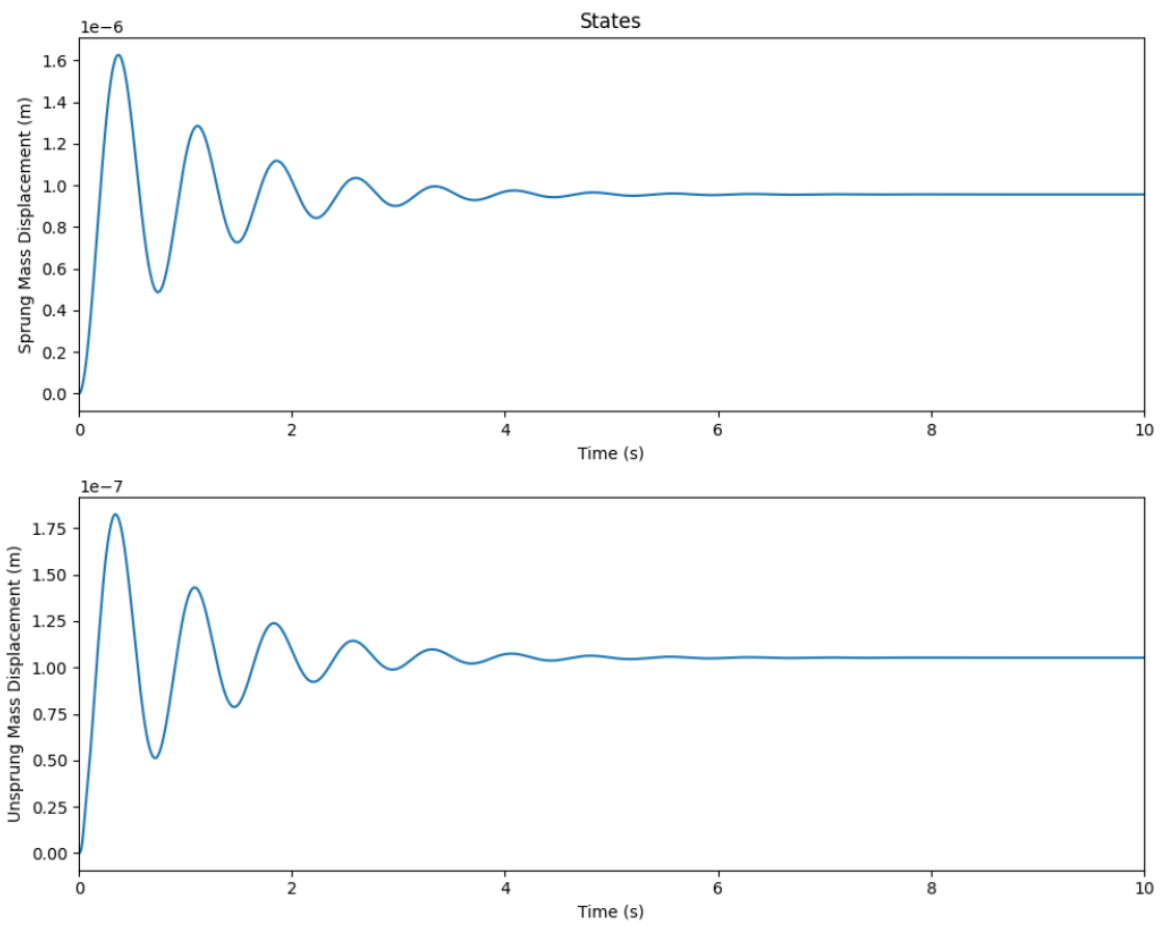
در این شبیه‌سازی، سیستم تعلیق به همراه ورودی‌های مرجع مانند موقعیت، سرعت و شتاب جاده مورد بررسی قرار می‌گیرد. ابتدا شرایط اولیه سیستم تنظیم شده و مقداردهی اولیه متغیرها انجام می‌شود. سپس در هر مرحله از شبیه‌سازی، بر اساس مقادیر ورودی، وضعیت جدید سیستم محاسبه شده و مقادیر حالت‌ها و خروجی ثبت می‌شوند. این فرآیند تا انتهای شبیه‌سازی ادامه می‌یابد و نتایج آن برای تحلیل رفتار سیستم ذخیره می‌شود.

در دو حالت مختلف، می‌توان شبیه‌سازی را اجرا کرد:

1. **حالت مرجع (Reference Mode)** که در آن خروجی مدل مرجع به‌عنوان مقدار ایده‌آل در نظر گرفته می‌شود.

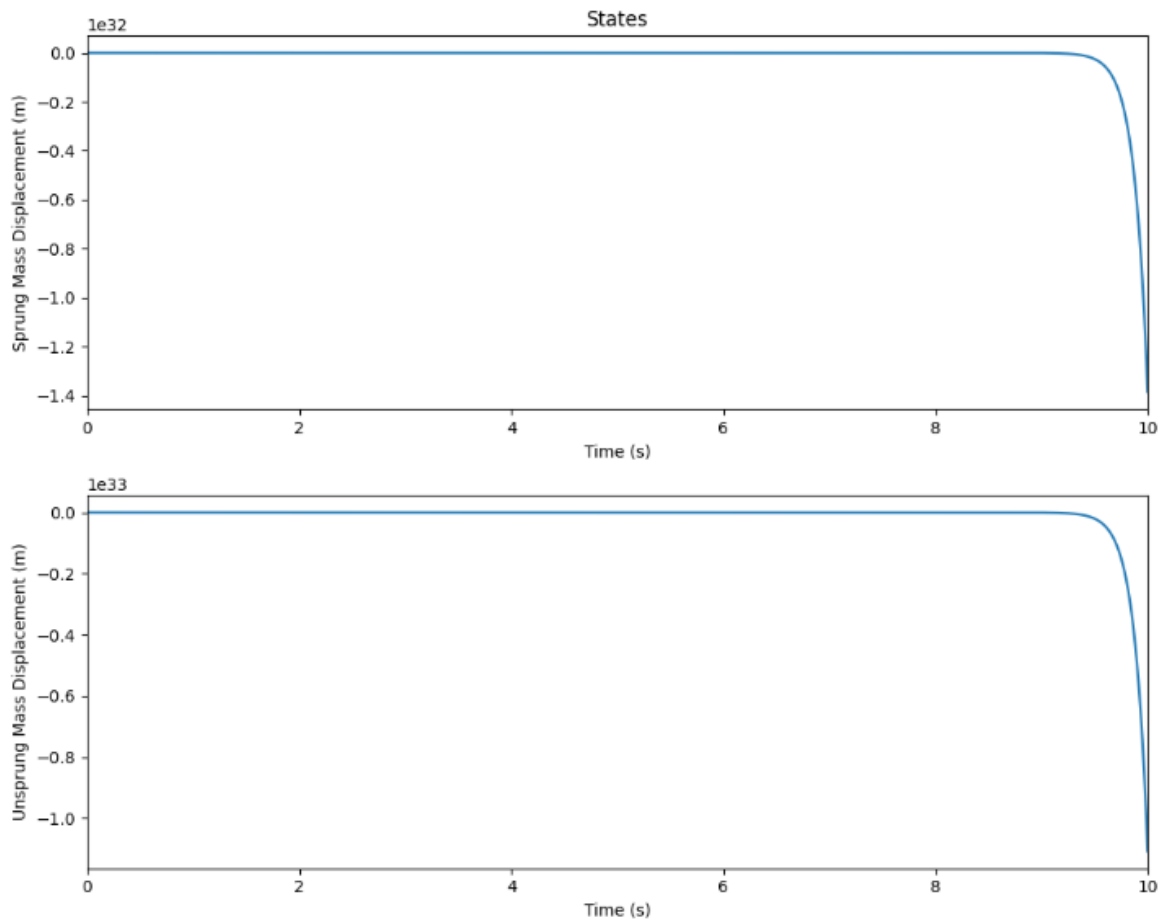
2. **حالت عادی** که در آن سیستم بدون در نظر گرفتن مدل مرجع تنها تحت تأثیر ورودی‌های داده‌شده واکنش نشان می‌دهد.

در نهایت، داده‌های شبیه‌سازی شامل **تاریخچه زمانی، مقادیر حالت‌ها و خروجی سیستم** در قالب یک دیکشنری جمع‌آوری شده و قابل استفاده برای تحلیل‌های بعدی می‌شوند. این نوع شبیه‌سازی، مبنای مقایسه عملکرد سیستم قبل و بعد از اعمال کنترل‌کننده را فراهم کرده و در طراحی سیستم‌های کنترلی مورد استفاده قرار می‌گیرد.

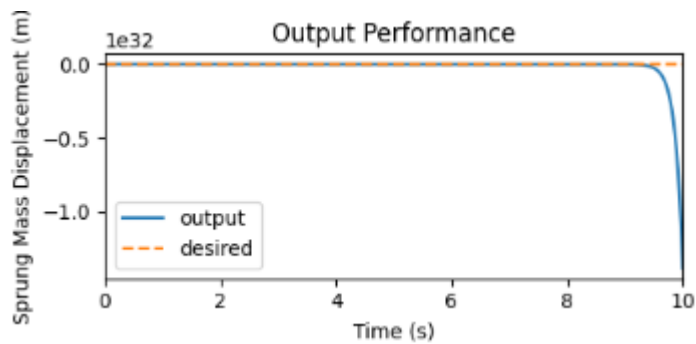


شکل 3-3 خروجی سیستم حلقه باز

همانطور که مشاهده میشود هم جسم متصل به فنر و هم جسم جدا از فنر میزان جابجایی میرا شونده‌ای دارند.

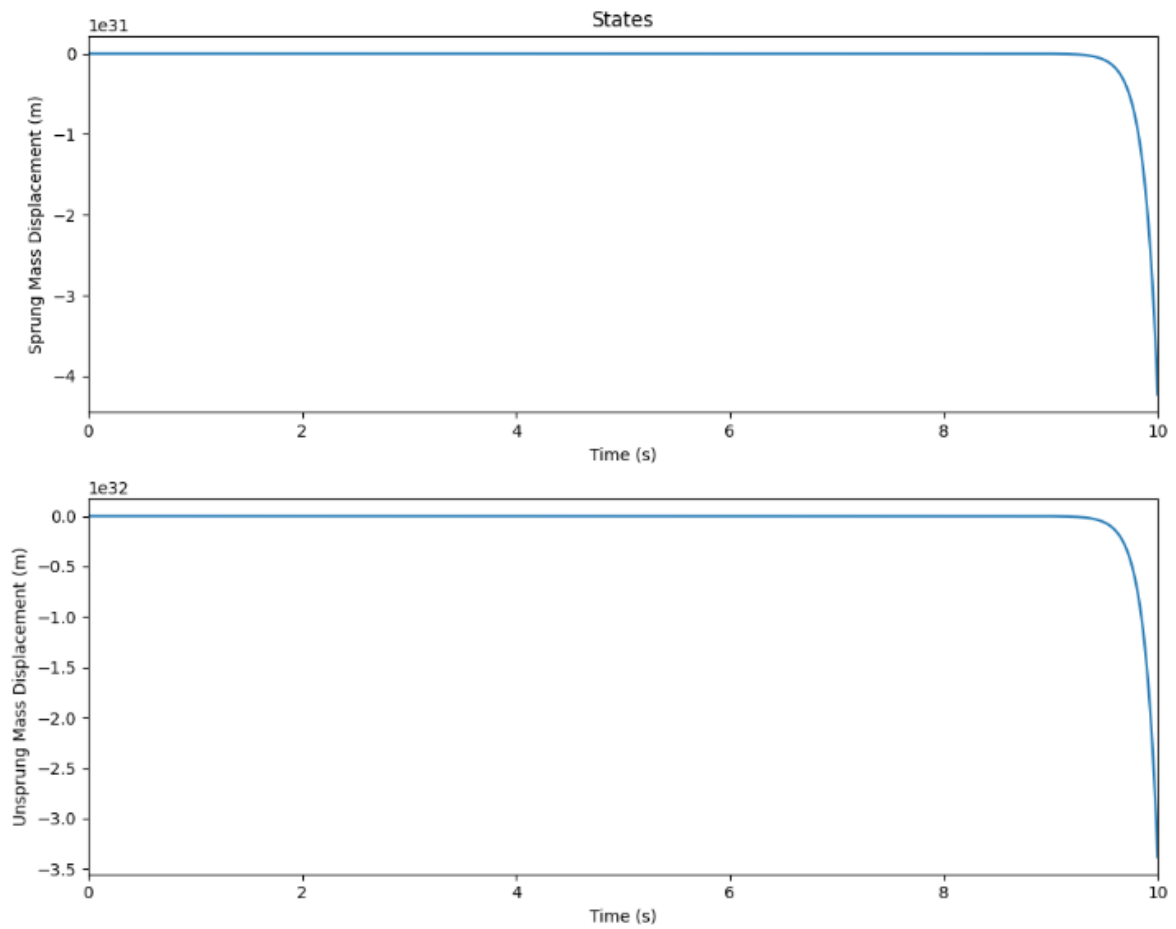


شکل 3-4 پاسخ حالات سیستم به ورودی پله

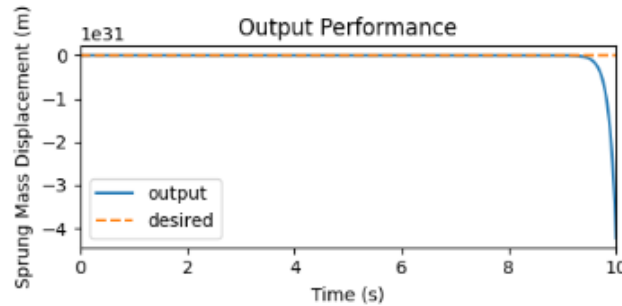


شکل 3-5 خروجی سیستم به ورودی پله

همانطور که مشاهده میکنید سیستم نسبت به ورودی پله ناپایدار می باشد.



شکل 3-5 پاسخ سیستم به ورودی سینوسی



شکل 3-5 خروجی سیستم نسبت به ورودی سینوسی

همانطور که مشاهده میشود سیستم نسبت به ورودی سینوسی ناپایدار است.

سیستم حلقه بسته

شبیه‌سازی حلقه بسته (**Closed-Loop Simulation**) یک روش ضروری در تحلیل و طراحی سیستم‌های کنترلی است که در آن سیستم، کنترل‌کننده و شناسایی‌کننده به صورت یکپارچه عمل می‌کنند. در این روش، خروجی سیستم به طور مداوم اندازه‌گیری شده و به کنترل‌کننده بازخورد داده می‌شود تا ورودی کنترلی مناسب را تولید کند. این فرآیند موجب بهبود عملکرد سیستم و کاهش خطا نسبت به مقدار مرجع می‌شود.

در این شبیه‌سازی، ابتدا مقادیر اولیه سیستم تنظیم شده و متغیرهای مورد نیاز مقداردهی اولیه می‌شوند. در هر گام زمانی، سیستم ورودی مرجع را دریافت کرده و کنترل‌کننده مبتنی بر شبکه عصبی **RBF** با توجه به مقدار خطای بین خروجی سیستم و مقدار مرجع، سیگنال کنترلی مناسب را تولید می‌کند. این سیگنال سپس به سیستم اعمال شده و باعث تغییر وضعیت آن می‌شود. علاوه بر این، یک شناسایی‌کننده **RBF** نیز در کنار سیستم قرار دارد که وظیفه یادگیری رفتار سیستم و بهبود عملکرد کنترل را بر عهده دارد.

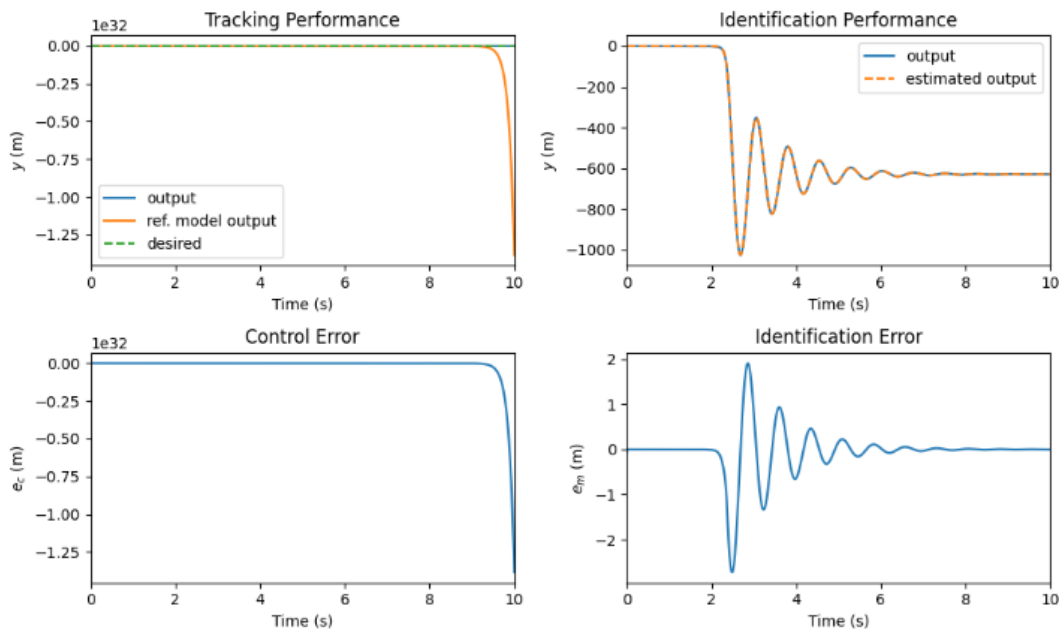
یکی از ویژگی‌های مهم این شبیه‌سازی، استفاده از ورودی‌های تأخیری برای کنترل و شناسایی است. برای این منظور، خروجی‌های گذشته و ورودی‌های کنترلی قبلی در قالب بردارهای تأخیری ذخیره

شده و به مدل کنترل کننده و شناسایی کننده داده می شوند. این تکنیک باعث بهبود دقت کنترل و پایداری سیستم می شود.

در هر گام، وزن های کنترل کننده و شناسایی کننده بر اساس خطاهای محاسبه شده به روزرسانی شده و سیستم به مرور عملکرد بهتری پیدا می کند. در پایان شبیه سازی، **تاریخچه زمانی، حالات سیستم، خروجی ها، مقادیر کنترلی و تغییرات وزن ها** ذخیره شده و برای تحلیل های بعدی قبل استفاده خواهد بود.

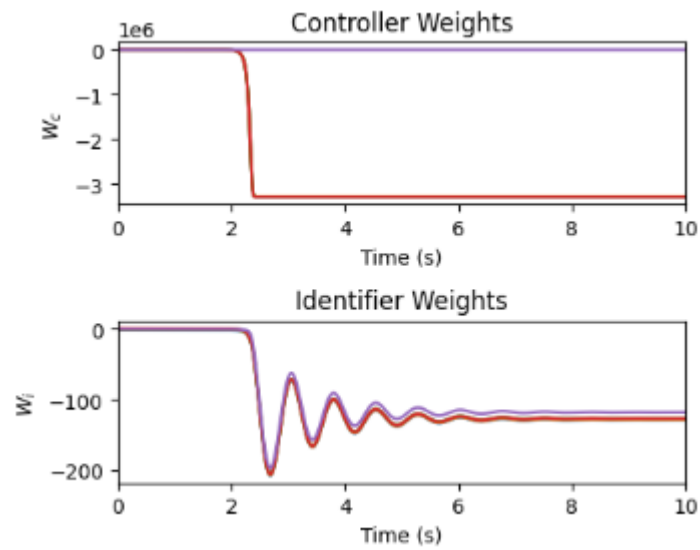
نتایج شبیه سازی

نتایج شبیه سازی با $\text{sampling_rate} = 0.001$ به شکل زیر می باشد.



شکل 3-6 نتایج شبیه سازی حلقه بسته

طبق خروجی ای که در بالا مشاهده میکنید شناسایی گر توانسته به خوبی سیستم را شناسایی کند و خطای آن میرا شونده می باشد اما کنترلر ناپایدار می باشد و نتیجه مطلوب را در tracking کسب نکردیم.

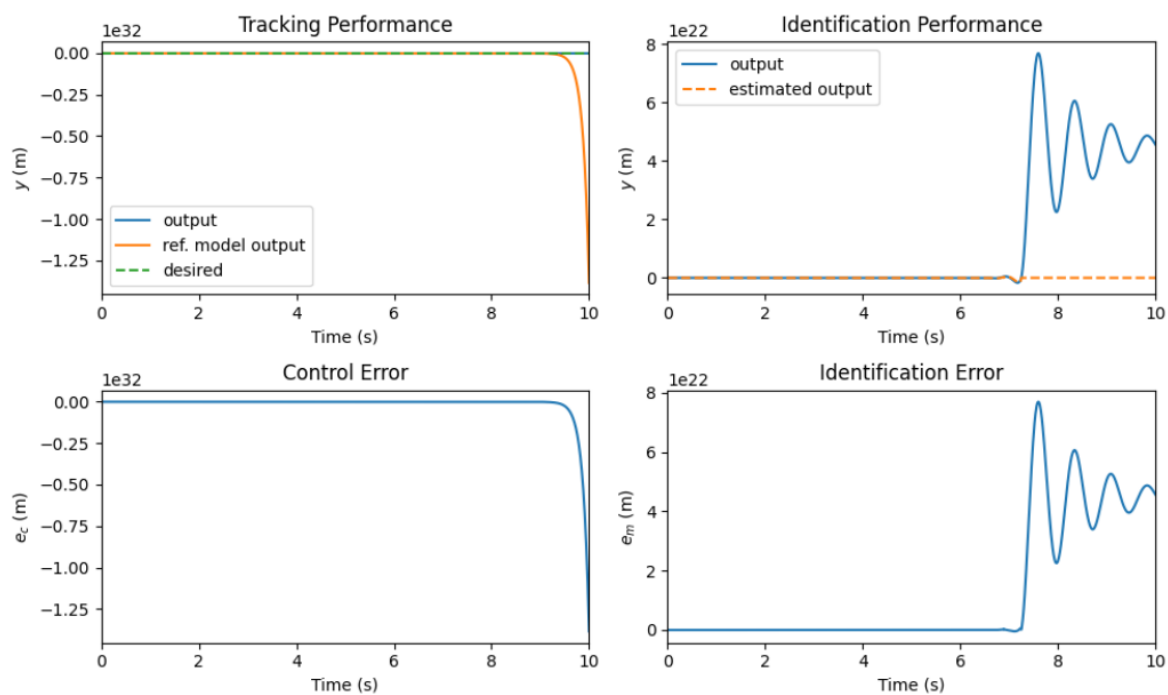


شکل 3-7 تغییرات وزن شناسایی کننده و کنترل کننده

همانطور که مشاهده میکنید وزن های شناسایی گر به صورت مطلوبی تغییر پیدا کرده اند اما از طرفی تغییرات وزن کنترلر نشان دهنده ناپایداری آن می باشد.

انجام آزمایش به ازای مقادیر مختلف sampling rate

$$T_d = 0.01$$



شکل 3-8 نتایج شبیه سازی با نرخ نمونه برداری 0.01 ثانیه

همانطور که مشاهده میکنید نتایج شناسایی گر همچنان میرا شونده می باشد اما در اینجا خروجی مطلوب نیست و کنترلر همچنان نا پایدار است.

لذا بهینه ترین میزان sampling rate را همان 0.001 در نظر می گیریم.

فصل چهارم:

کنترل مبتنی بر شبکه LSTM

معرفی شبکه LSTM

در این پروژه، یک مدل کنترلی مبتنی بر LSTM برای یک سیستم تعلیق چهار درجه آزادی توسعه داده شده است. هدف از این تحقیق، طراحی یک کنترل کننده هوشمند است که بتواند عملکرد سیستم تعلیق را بهبود بخشیده و نوسانات ناشی از ناهمواری‌های جاده را به حداقل برساند. مدل پیشنهادی از داده‌های ورودی شامل نیروهای اعمالی و پاسخ‌های سیستم استفاده کرده و از طریق یادگیری توالی‌های زمانی، یک استراتژی کنترلی بهینه را ارائه می‌دهد.

در ادامه، ابتدا به معرفی مدل دینامیکی سیستم تعلیق پرداخته شده، سپس روش طراحی و آموزش شبکه LSTM تشریح می‌شود. در نهایت، عملکرد کنترل کننده پیشنهادی با استفاده از شبیه‌سازی‌های عددی ارزیابی می‌گردد.

این نکته حائز اهمیت می باشد که برای دریافت نتایج بهتر از 4 درجه آزادی راننده در مدل صرف نظر شده است و صرفاً 4 درجه آزادی مربوط به خودرو در نظر گرفته شده است.

ساختار شبکه LSTM

LSTM یک نوع شبکه عصبی بازگشتی (RNN) است که برای مدل‌سازی داده‌های دنباله‌دار و وابسته به زمان طراحی شده است. برخلاف RNN های معمولی که با مشکل ناپدید شدن گرادیان مواجه هستند، LSTM دارای سلول‌های حافظه و دروازه‌های کنترلی است که امکان یادگیری روابط طولانی مدت را فراهم می‌کنند.

ساختار LSTM

در هر واحد LSTM، چهار لایه اصلی وجود دارد که ورودی و وضعیت قبلی را پردازش می‌کنند:

1. دروازه فراموشی (Forget Gate)

- تصمیم می‌گیرد که چه اطلاعاتی از حالت قبلی حذف شود.

- با استفاده از یک تابع سیگموید، مقدار بین ۰ تا ۱ را تولید کرده و تعیین می‌کند که چه مقدار از اطلاعات قبلی باید نگه داشته شود.

2. دروازه ورودی (Input Gate)

- مشخص می‌کند که چه اطلاعات جدیدی به حافظه اضافه شود.
- از دو بخش تشکیل شده: یک لایه سیگموید که تصمیم می‌گیرد چه بخش‌هایی از داده‌های جدید مهم هستند، و یک لایه تانژانت هایپربولیک که مقدار جدید را ایجاد می‌کند.

3. به‌روزرسانی حافظه سلولی

- با ترکیب اطلاعات جدید و اطلاعات قبلی، مقدار جدیدی برای حافظه سلولی محاسبه می‌شود.

4. دروازه خروجی (Output Gate)

- مشخص می‌کند که چه مقداری از حالت سلولی به خروجی ارسال شود.
- از ترکیب لایه سیگموید و تابع تانژانت هایپربولیک برای محاسبه خروجی استفاده می‌کند.

نمایش گرافیکی ساختار LSTM

در تصویر زیر، معماری یک واحد LSTM نمایش داده شده است:

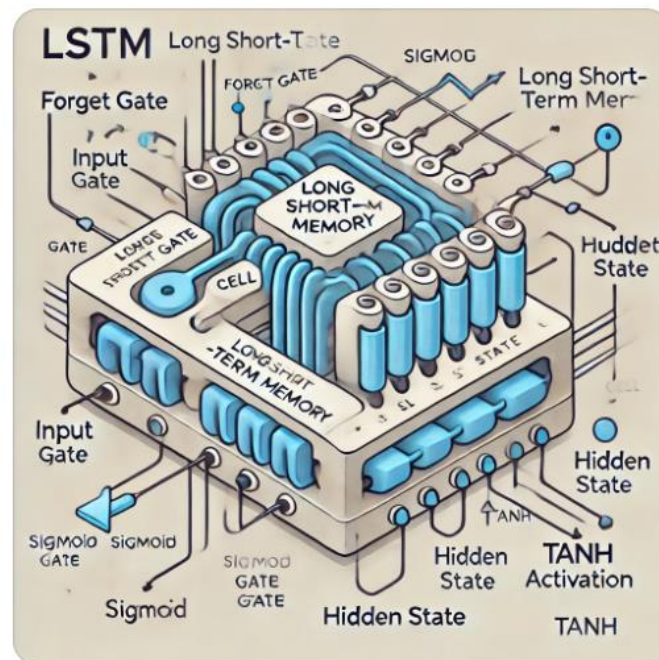
C: بردار حافظه سلولی

h: بردار خروجی

x: ورودی در هر زمان

σ : تابع سیگموید

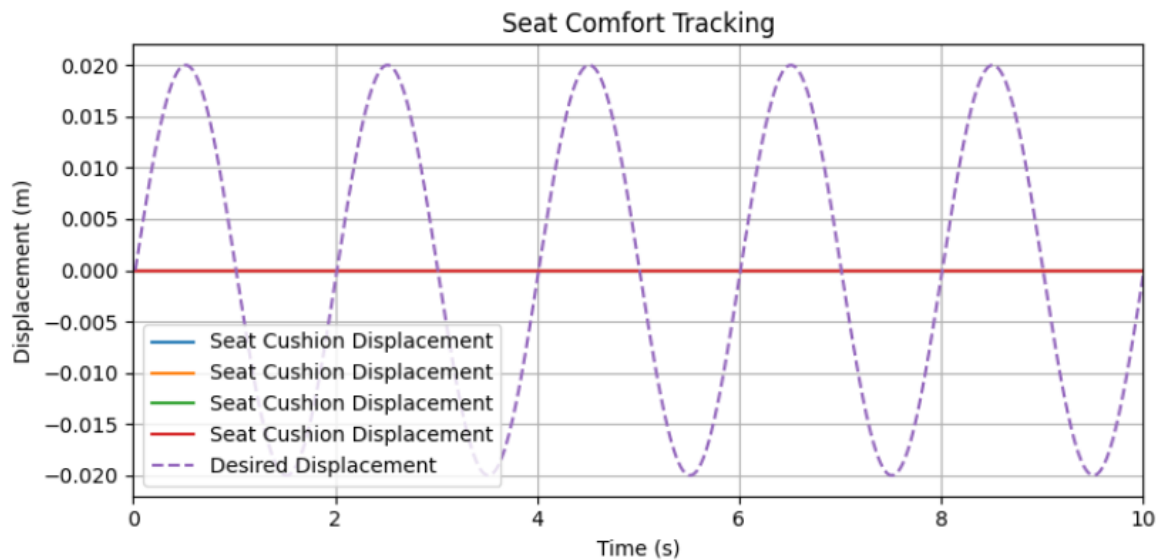
\tanh : تابع تانژانت هایپربولیک



شکل 1-4 ساختار شبکه عصبی LSTM

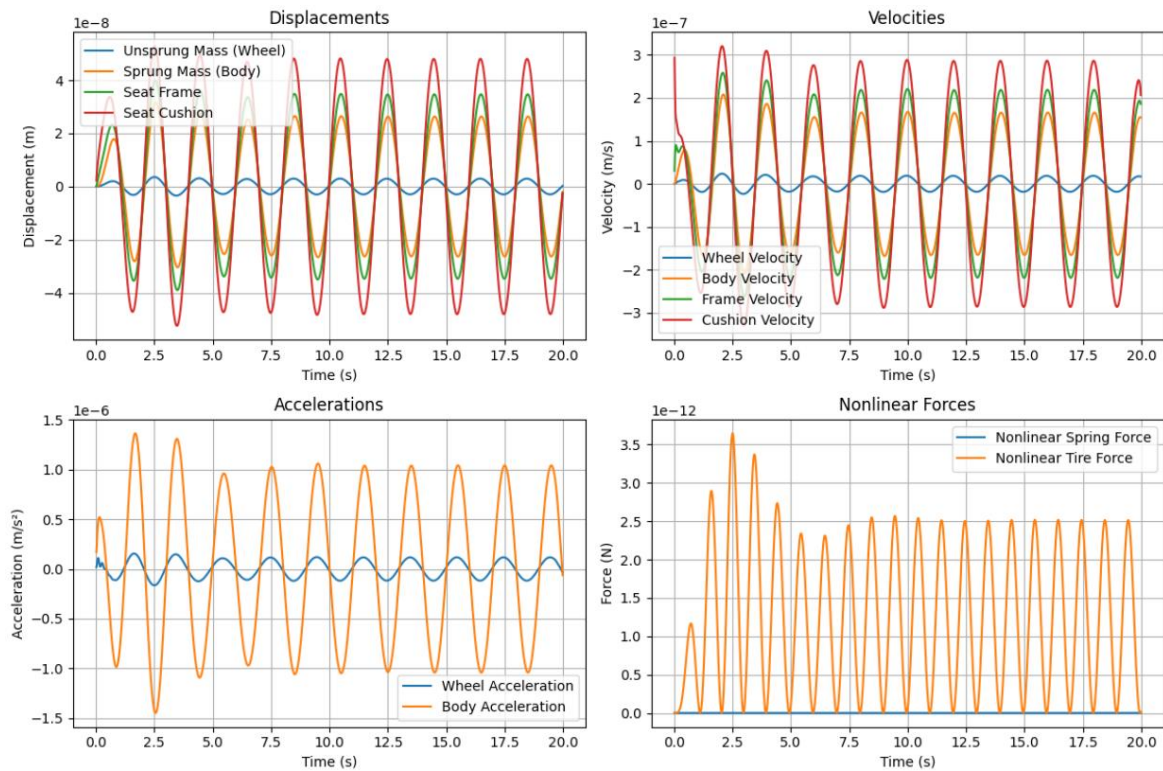
طراحی مدل و پیاده سازی مدل غیر خطی

در این بخش با استفاده از معادلات دینامیکی سیستم که در مقاله موجود می باشد ، سیستم را مدلسازی کردیم



شکل 2-4 مدلسازی سیستم غیر خطی

همانطور که مشاهده می کنید جابجایی صندلی نسبت به حالت مطلوب نا مناسب بوده و نیازمند این می باشد تا با پیاده سازی کنترلر و رویتگر مناسب تعقیب به خوبی انجام شود.



شکل 3-4 نتایج مدلسازی سیستم (حالات سیستم)

با توجه به تصویر بالا می توان گفت که مدلسازی به نحو خوبی انجام شده است.

این تصویر مجموعه ای از نمودارهای شبیه سازی دینامیکی یک سیستم تعلیق خودرو را نشان می دهد که قبل از پیاده سازی کنترلر و رویت گر انجام شده است. این سیستم شامل جرم های فنربندی شده (بدنه)، جرم های فنربندی نشده (چرخ)، قاب صندلی و بالشک صندلی است. چهار نمودار مجزا برای نمایش تغییرات جابجایی، سرعت، شتاب و نیروهای غیرخطی در طول زمان ارائه شده اند.

تحلیل نمودارها:

1. نمودار بالا-چپ (Displacements)

- نشان دهنده میزان جابجایی های اجزای مختلف سیستم تعلیق در واحد متر است.

- مشاهده می‌شود که بالشتک صندلی (Seat Cushion) بیشترین میزان نوسان را دارد، در حالی که چرخ (Unsprung Mass) دارای کمترین دامنه جابجایی است.
- این نمودار نشان‌دهنده رفتار کلی سیستم تعلیق در برابر تحریک ورودی (مانند ناهمواری‌های جاده) است.

2. نمودار بالا-راست (Velocities)

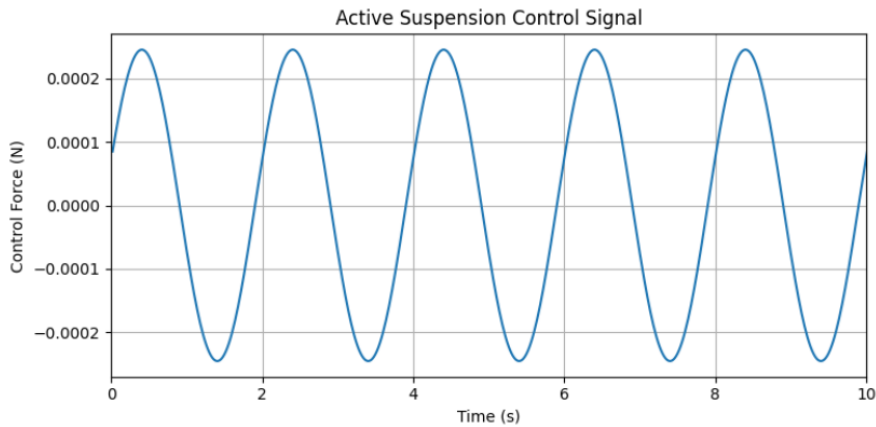
- سرعت‌های اجزای مختلف سیستم نمایش داده شده‌اند.
- مشاهده می‌شود که بیشترین دامنه سرعت مربوط به بالشتک صندلی و کمترین آن مربوط به چرخ است.
- شکل موج‌های مشابه نمودار جابجایی‌ها اما با تغییر فاز مشاهده می‌شود.

3. نمودار پایین-چپ (Accelerations)

- شتاب بدنه (Body Acceleration) و شتاب چرخ (Wheel Acceleration) نمایش داده شده‌اند.
- مشاهده می‌شود که شتاب بدنه نسبت به شتاب چرخ دامنه بیشتری دارد، که نشان‌دهنده تأثیر مستقیم نوسانات جاده بر روی بدنه خودرو است.
- نوسانات شتاب، عامل تأثیرگذار بر راحتی سرنشینان خودرو هستند.

4. نمودار پایین-راست (Nonlinear Forces)

- نمایش‌دهنده نیروهای غیرخطی موجود در سیستم، از جمله نیروی فنر غیرخطی و نیروی تایر غیرخطی است.
- مشاهده می‌شود که نیروی تایر غیرخطی دارای نوسانات شدیدتری نسبت به نیروی فنر غیرخطی است.
- این نیروها به دلیل وجود المان‌های غیرخطی مانند لاستیک و فنرهای سیستم تعلیق ایجاد می‌شوند.



شکل 4-4 نتایج مدل‌سازی سیستم (سیگنال کنترلی)

این تصویر نمودار سیگنال کنترلی سیستم تعلیق فعال را نمایش می‌دهد. محور افقی نشان‌دهنده زمان (Time) بر حسب ثانیه است و محور عمودی مقدار نیروی کنترلی (Control Force) را بر حسب نیوتن نمایش می‌دهد.

تحلیل نمودار:

1. ماهیت نوسانی سیگنال کنترلی:

- نیروی کنترلی دارای یک شکل موج سینوسی متناوب است که نشان‌دهنده اعمال یک نیروی متغیر برای تنظیم عملکرد سیستم تعلیق می‌باشد.
- این نوع سیگنال معمولاً در سیستم‌های کنترلی که از کنترلرهای بازخوردی مانند کنترلر **LQR (Linear Quadratic Regulator)** یا **PID** استفاده می‌کنند، دیده می‌شود.

2. دامنه نیروی کنترلی:

- مقدار نیروی کنترلی در بازه‌ای بین تقریباً -0.0002 تا 0.0002 نیوتن تغییر می‌کند.

- این مقدار نشان می‌دهد که کنترلر نیروی محدودی را برای کاهش نوسانات سیستم تعلیق اعمال می‌کند.

3. دوره تناوب و فرکانس سیگنال:

- دوره تناوب (زمان یک سیکل کامل نوسان) حدود 2 ثانیه است که نشان‌دهنده فرکانس پایین سیگنال کنترلی است.
- این امر معمولاً نشان‌دهنده تنظیمات کنترلی است که هدف آن بهبود پایداری و راحتی سرنشینان بدون اعمال نیروهای شدید است.

تولید مجموعه داده برای آموزش مدل در سیستم تعلیق فعال خودرو

روش تولید داده

در این روش، داده‌های مورد نیاز با اجرای شبیه‌سازی‌های متعددی از سیستم تعلیق خودرو تولید می‌شوند. شبیه‌سازی‌ها با شرایط اولیه تصادفی انجام شده و ورودی‌های کنترلی مشخصی به سیستم اعمال می‌شود. سپس، خروجی‌های سیستم و وضعیت‌های داخلی آن در طول زمان ثبت شده و به عنوان مجموعه داده مورد استفاده قرار می‌گیرند.

ساختار مجموعه داده

مجموعه داده تولید شده شامل سه بخش اصلی است:

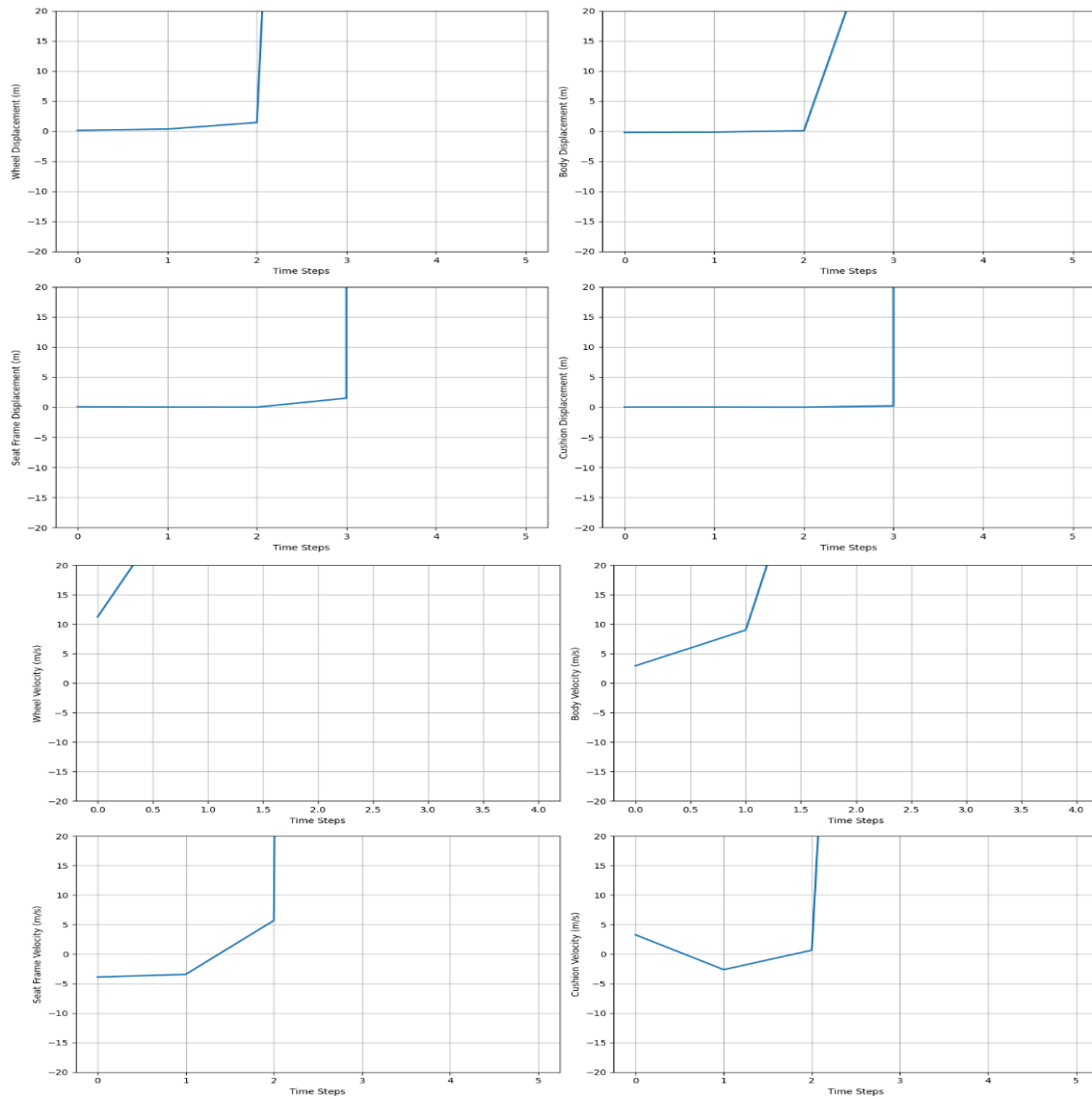
1. **ورودی‌های کنترلی:** نیروی اعمال شده توسط سیستم تعلیق فعال که به دو صورت **جاده‌ای (Road Disturbance)** و **پله‌ای (Step Input)** تولید می‌شود.
2. **وضعیت‌های داخلی:** مقادیر مربوط به متغیرهای داخلی سیستم تعلیق که شامل جابجایی‌ها، سرعت‌ها و شتاب‌های بخش‌های مختلف سیستم است.
3. **خروجی‌های اندازه‌گیری شده:** داده‌هایی که از حسگرهای سیستم جمع‌آوری می‌شوند و نشان‌دهنده رفتار کلی سیستم در پاسخ به ورودی‌های کنترلی هستند.

پیش‌پردازش داده‌ها

به منظور بهبود کیفیت داده‌ها و افزایش دقت مدل‌های یادگیری ماشین، داده‌های تولید شده نرمال‌سازی می‌شوند. این فرآیند باعث می‌شود که متغیرهای ورودی و خروجی در یک بازه مشخص مقیاس‌بندی شده و تأثیر مقادیر بسیار بزرگ یا بسیار کوچک بر آموزش مدل کاهش یابد. علاوه بر این، داده‌ها به گونه‌ای سازماندهی می‌شوند که بتوان از آن‌ها برای دو هدف استفاده کرد:

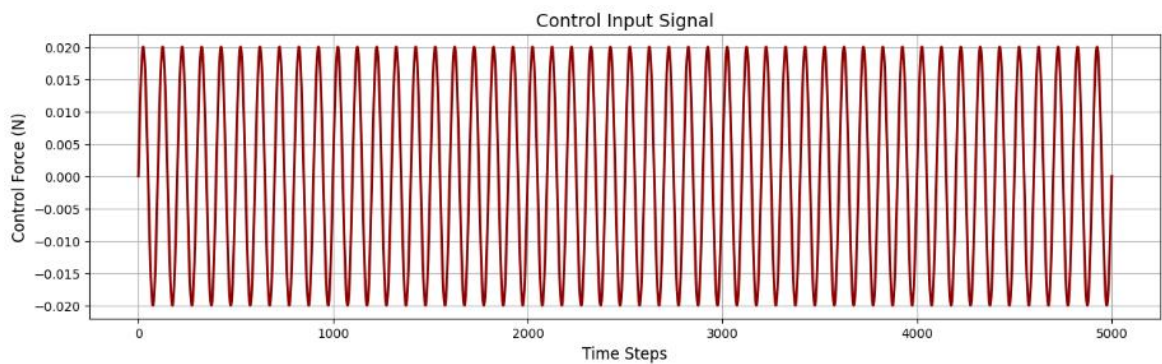
- **مدل‌سازی دینامیکی:** پیش‌بینی وضعیت آینده سیستم بر اساس وضعیت‌های قبلی و ورودی‌های کنترلی.

- طراحی روییت گر: تخمین وضعیت‌های داخلی سیستم بر اساس خروجی‌های اندازه‌گیری شده و ورودی‌های کنترلی.



شکل 4-5 خروجی دیتا سازی

پس از تولید مجموعه داده‌ها، نتایج شبیه‌سازی برای بررسی پاسخ سیستم تعلیق خودرو تحلیل شده‌اند. نمودارهای ارائه‌شده تغییرات جابه‌جایی و سرعت اجزای مختلف سیستم، از جمله چرخ، بدنه، فریم صندلی و کوسن را در طول زمان نشان می‌دهند. رفتار سیستم در واکنش به ورودی‌های کنترلی و تحریکات جاده‌ای بررسی شده و ناپیوستگی‌های موجود در داده‌ها نشان‌دهنده تغییرات ناگهانی ناشی از ورودی‌های کنترلی است. این تحلیل‌ها برای ارزیابی دقت مدل‌سازی و عملکرد سیستم تعلیق در شرایط مختلف ضروری هستند.



شکل 4-6 ورودی کنترلی

این نمودار سیگنال ورودی کنترلی را نشان می‌دهد که به سیستم تعلیق اعمال شده است. این سیگنال یک موج متناوب با دامنه محدود در محدوده $[-0.02, 0.02]$ نیوتن است که با یک فرکانس ثابت تغییر می‌کند. هدف از اعمال این سیگنال بررسی واکنش سیستم تعلیق به تحریک‌های دوره‌ای و ارزیابی عملکرد کنترل‌کننده در کاهش ارتعاشات و بهبود پایداری خودرو می‌باشد.

همانطور که مشاهده می کنید مدل به شدت ناپایدار بوده و حتی با حذف 4 درجه آزادی مربوط به راننده خروجی نمودارهای حالات بشدت نا پایدار می باشد و در آزمایش اول سیستم موفق به تعقیب سیگنال مورد انتظار نشده است.

این نتیجه نشان می دهد مدل سازی سیستم با این روش تقریبا غیر ممکن می باشد.

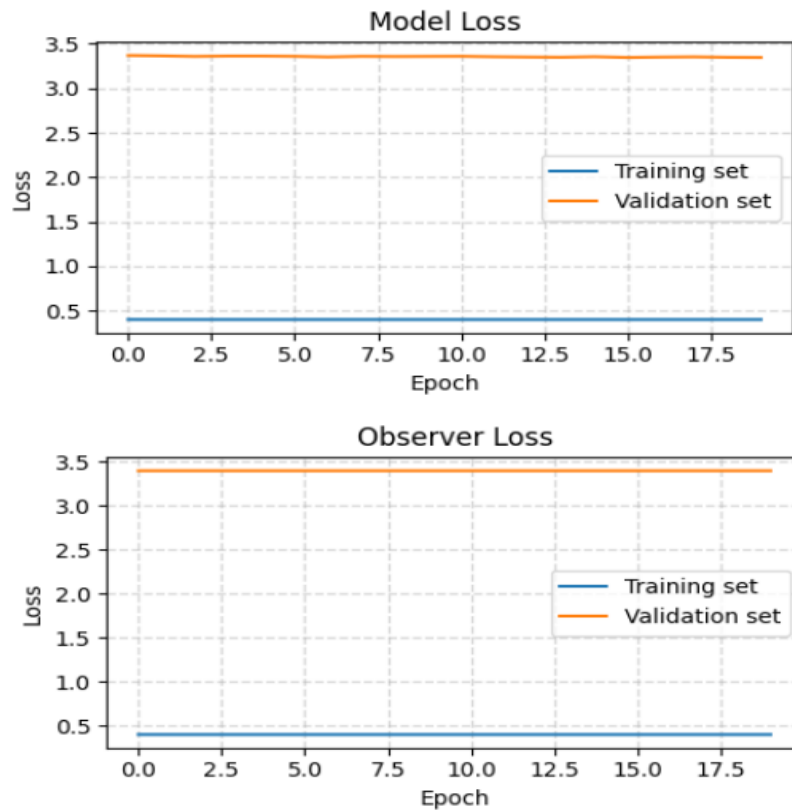
مدل پیشنهادی برای شناسایی و رویت سیستم

در این پژوهش، از شبکه حافظه دار بلندمدت (LSTM) به عنوان یک مدل داده محور برای شناسایی سیستم و طراحی رویتگر استفاده شده است. دو مدل مجزا توسعه داده شده اند: یک مدل LSTM برای شناسایی رفتار دینامیکی سیستم و یک مدل دیگر برای طراحی رویتگر جهت تخمین متغیرهای غیرقابل مشاهده.

مدل شناسایی سیستم بر اساس ورودی های مشاهده شده و خروجی های متناظر، ساختاری عمیق شامل یک لایه LSTM و لایه های متراکم (Dense) برای استخراج ویژگی های وابسته به زمان دارد. این مدل به منظور پیش بینی حالات آینده سیستم با استفاده از داده های آموزشی مورد استفاده قرار می گیرد.

در مقابل، رویتگر مبتنی بر LSTM وظیفه تخمین حالات سیستم را بر اساس ورودی های مشاهده شده بر عهده دارد. این مدل با ساختاری مشابه مدل شناسایی، اما با معماری مناسب برای تخمین حالات پنهان، امکان بازسازی متغیرهای دینامیکی سیستم را فراهم می کند.

برای بهبود دقت یادگیری، داده ها قبل از ورود به شبکه نرمال سازی شده اند و از تکنیک هایی مانند تنظیم نرخ یادگیری و کاهش نوسانات گرادیان در بهینه سازی استفاده شده است. نتایج اولیه نشان می دهد که این روش می تواند در مدل سازی سیستم های دینامیکی عملکرد قابل قبولی ارائه دهد، اما همچنان محدودیت هایی در مقابله با ناپایداری ها و پیچیدگی های سیستم تعلیق پیشرفته مشاهده شده است.



شکل 4-7 نتایج آموزش مدل

نتایج آموزش مدل LSTM برای شناسایی و رویت سیستم نشان‌دهنده عملکرد نامناسب این مدل در تخمین پارامترهای سیستم است. همان‌طور که در نمودارهای **Model Loss** و **Observer Loss** مشاهده می‌شود، مقدار **خطای اعتبارسنجی** (Validation Loss) در تمامی دوره‌های آموزشی تقریباً ثابت و در سطح بالایی باقی مانده است. این موضوع بیانگر آن است که مدل نتوانسته است الگوی مناسبی از داده‌های آموزشی استخراج کند و در نتیجه، تعمیم‌پذیری آن برای داده‌های جدید بسیار ضعیف خواهد بود.

از سوی دیگر، اختلاف قابل توجه میان **خطای آموزش** (Training Loss) و **خطای اعتبارسنجی** بیانگر آن است که مدل دچار **overfitting** نشده، بلکه اساساً در یادگیری الگوهای سیستم ناتوان بوده است. این ضعف نشان می‌دهد که معماری انتخاب‌شده، میزان داده‌ها، و یا فرآیند تنظیم مدل، نتوانسته‌اند مدل مناسبی برای سیستم مورد نظر ارائه دهند.

بنابراین، استفاده از این مدل برای کاربردهای کنترلی یا شناسایی سیستم نامناسب بوده و نیاز به روش‌های جایگزین مانند **مدل‌های کلاسیک کنترل**، **تخمینگرهای مبتنی بر فیلتر کالمن**، یا **روش‌های یادگیری عمیق پیشرفته‌تر** (مانند شبکه‌های Transformer یا Hybrid Models) احساس می‌شود.

نتیجه گیری کلی

در این پژوهش، از شبکه‌های بازگشتی LSTM برای شناسایی و رویت سیستم تعلیق پیشرفته استفاده شد. هدف از این مطالعه، بررسی کارایی این روش در تخمین دینامیک سیستم و طراحی یک رویکرد کارآمد برای تخمین متغیرهای حالت بود. با این حال، نتایج حاصل از آموزش مدل نشان داد که روش پیشنهادی دارای ضعف‌های جدی بوده و نمی‌تواند عملکرد مطلوبی را در این زمینه ارائه دهد.

1. عملکرد ضعیف در فرآیند آموزش

نمودارهای مربوط به میزان خطای مدل و رویکرد نشان داد که مقدار خطای اعتبارسنجی (Validation Loss) در طول کل دوره‌های آموزشی تقریباً ثابت باقی مانده و در سطح بالایی قرار دارد. این پدیده نشان‌دهنده آن است که مدل نتوانسته است الگوی مناسب و معناداری را از داده‌های آموزشی استخراج کند. به عبارت دیگر، فرآیند یادگیری مدل ناقص بوده و شبکه از لحاظ تعمیم‌پذیری ناتوان است.

2. ضعف در تعمیم‌پذیری و همگرایی مدل

یکی از مشکلات اساسی مشاهده‌شده در نتایج، اختلاف ناچیز بین میزان خطای آموزش (Training Loss) و اعتبارسنجی بود که معمولاً بیانگر عدم یادگیری صحیح شبکه است. در یک مدل بهینه، انتظار می‌رود که با افزایش تعداد دوره‌های آموزشی، میزان خطای مدل کاهش یافته و همگرایی در یادگیری اتفاق بیفتد. اما در اینجا مدل در مرحله آموزش حتی به یک حداقل موضعی مناسب نیز نرسیده است. این نشان می‌دهد که شبکه LSTM برای مدل‌سازی دینامیک سیستم تعلیق پیشرفته مناسب نبوده و نیازمند معماری‌های جایگزین است.

3. ناکارآمدی در تخمین متغیرهای سیستم

مدل طراحی شده نه تنها در پیش‌بینی متغیرهای خروجی سیستم دچار مشکل است، بلکه رویکرد مبتنی بر LSTM نیز قادر به ارائه تخمین مناسبی از متغیرهای حالت سیستم نیست. این مسئله در عمل موجب عدم اطمینان در پیاده‌سازی واقعی می‌شود، چرا که رویکردها معمولاً باید بتوانند نویزها و ناپایداری‌های محیطی را نیز در نظر گرفته و عملکردی پایدار ارائه دهند.

راهکارهای پیشنهادی

با توجه به نتایج به دست آمده، پیشنهاد می شود که به جای استفاده از **LSTM**، روش های جایگزین مورد بررسی قرار گیرند. برخی از گزینه های پیشنهادی عبارت اند از:

- روش های کلاسیک کنترل و فیلترگذاری: استفاده از فیلتر کالمن توسعه یافته (EKF) یا فیلتر کالمن غیرخطی (UKF) برای بهبود رویت سیستم های پیچیده. این روش ها برای سیستم های دینامیکی مانند سیستم تعلیق عملکرد بهتری نسبت به مدل های یادگیری عمیق خواهند داشت.
- شبکه های ترکیبی: ترکیب مدل های کلاسیک کنترلی با یادگیری عمیق، مانند استفاده از شبکه های عصبی ترکیبی (Hybrid Models) که داده های استخراج شده از مدل های فیزیکی را بهبود می بخشند.
- روش های یادگیری عمیق پیشرفته: استفاده از مدل هایی نظیر **Transformer** یا **Attention-based Networks** که در پردازش سری های زمانی عملکرد بهتری نسبت به LSTM دارند.
- بهینه سازی معماری مدل: بررسی افزایش تعداد لایه ها، تغییر در توابع فعال سازی، استفاده از مکانیزم های توجه (Attention Mechanisms) و تنظیم دقیق تر هایپر پارامترها برای بهبود فرآیند یادگیری.

فصل پنجم:

کنترل مبتنی بر RL

معرفی روش کنترل مبتنی بر RL

پروژه کنترل هوشمند سیستم‌های تعلیق خودرو یکی از چالش‌های پیچیده در مهندسی خودرو است که می‌تواند تاثیر زیادی بر راحتی راننده و پایداری خودرو داشته باشد. یکی از روش‌های نوین و پیشرفته در این زمینه استفاده از الگوریتم‌های یادگیری تقویتی (Reinforcement Learning) یا (RL) است. این الگوریتم‌ها با امکان یادگیری از محیط و بهینه‌سازی مستمر عملکرد سیستم، می‌توانند راه‌حل‌های بهینه‌تری برای کنترل سیستم‌های پیچیده مانند سیستم تعلیق خودرو ارائه دهند.

در این گزارش، به تحلیل و طراحی یک سیستم کنترل هوشمند برای یک مدل سیستم تعلیق 4 درجه آزادی خودرو پرداخته شده است. برای این منظور، از الگوریتم Soft Actor-Critic (SAC) به عنوان یک روش پیشرفته یادگیری تقویتی استفاده شده است. این الگوریتم به دلیل توانایی بالای آن در یادگیری سیاست‌های بهینه در محیط‌های پیوسته و پیچیده، انتخاب مناسبی برای کنترل سیستم تعلیق است. سیستم تعلیق 4 درجه آزادی به دلیل پیچیدگی‌های مدل و نیاز به دقت بالا در تنظیمات، چالش‌های زیادی را در کنترل به همراه دارد که استفاده از متدهای هوشمند مانند SAC می‌تواند به بهبود عملکرد آن کمک کند.

در این گزارش، ابتدا به معرفی سیستم تعلیق 4 درجه آزادی و ویژگی‌های آن پرداخته می‌شود، سپس به بررسی الگوریتم SAC و نحوه پیاده‌سازی آن برای کنترل سیستم تعلیق خودرو خواهیم پرداخت. هدف از این تحقیق بهبود کارایی سیستم تعلیق از طریق یادگیری ماشین و الگوریتم‌های پیشرفته به منظور ارتقاء راحتی و ایمنی خودرو است.

معرفی الگوریتم SAC

الگوریتم Soft Actor-Critic (SAC) یک الگوریتم پیشرفته یادگیری تقویتی است که به طور خاص برای محیط‌های پیوسته طراحی شده است. این الگوریتم به عنوان یک الگوریتم actor-critic شناخته می‌شود، که در آن دو شبکه عصبی مجزا برای مدل‌سازی سیاست (actor) و تابع ارزش (critic) استفاده می‌شود. یکی از ویژگی‌های برجسته SAC، استفاده از استراتژی‌های نرم (soft) برای بهبود پایداری و بهره‌وری در یادگیری است.

SAC برخلاف الگوریتم‌های کلاسیک که از سیاست‌های قطعی (deterministic) استفاده می‌کنند، از سیاست‌های تصادفی (stochastic) بهره می‌برد. این به این معناست که در هر گام از فرآیند یادگیری، سیاست انتخابی به صورت تصادفی از توزیع‌های احتمالی انتخاب می‌شود که به طور موثر باعث افزایش تنوع و توانایی کشف محیط می‌شود.

یکی دیگر از ویژگی‌های کلیدی SAC، استفاده از مفهوم "entropy" در فرآیند یادگیری است. این مفهوم باعث می‌شود که الگوریتم به طور همزمان نه تنها به حداکثر کردن پاداش‌ها بلکه به حداقل کردن پیش‌بینی‌پذیری سیاست (یعنی افزایش تصادفی بودن آن) نیز توجه داشته باشد. این رویکرد کمک می‌کند تا الگوریتم بتواند بهتر به اکتشاف محیط بپردازد و از بیش‌برازش (overfitting) جلوگیری کند. SAC به طور گسترده در مسائل پیچیده با محیط‌های پیوسته و غیرخطی مورد استفاده قرار می‌گیرد، زیرا توانایی یادگیری بهینه سیاست‌ها و بهینه‌سازی عملکرد سیستم‌ها را در چنین محیط‌هایی به خوبی دارد.

برای اطلاعات بیشتر از روابط به نوت بوک نهایی که به پیوست فایل ارسال شده است مراجعه کنید.

توضیح بخش های مختلف برنامه

توضیح عملکرد کلاس ReplayBuffer

کلاس **ReplayBuffer** وظیفه ذخیره سازی و مدیریت تجربیات محیط را بر عهده دارد. در فرآیند یادگیری تقویتی، این تجربیات شامل وضعیت های مختلف سیستم، اقدامات انجام شده، پاداش های دریافتی، وضعیت های بعدی و اطلاعاتی در مورد پایان تعامل است. این بافر به الگوریتم کمک می کند تا از تجربیات گذشته برای بهبود سیاست یادگیری استفاده کند. عملکرد کلی این کلاس به شرح زیر است:

1. ذخیره تجربیات

این بخش به طور مداوم تجربیات جدید را ذخیره می کند. هر تجربه شامل وضعیت فعلی، اقدام انجام شده، پاداش دریافت شده، وضعیت بعدی و مشخصه پایان تعامل است. این اطلاعات به طور مرتب در یک فضای ذخیره سازی ذخیره می شود تا در مراحل بعدی برای یادگیری از آن ها استفاده شود.

2. نمونه برداری تصادفی

در این بخش، تجربیات ذخیره شده به طور تصادفی انتخاب می شوند تا در فرآیند یادگیری استفاده شوند. انتخاب تصادفی تجربیات به الگوریتم این امکان را می دهد که یادگیری را بدون وابستگی به ترتیب تجربیات انجام دهد و از همبستگی میان نمونه ها جلوگیری کند.

3. محدود کردن اندازه بافر

برای جلوگیری از مصرف بیش از حد حافظه، تعداد تجربیات ذخیره شده در بافر محدود به یک مقدار مشخص است. وقتی تعداد تجربیات به حد مجاز برسد، تجربیات قدیمی تر جایگزین تجربیات جدیدتر می شوند.

4. میزان تجربیات ذخیره شده

این بخش تعداد تجربیات ذخیره شده در بافر را برمی گرداند، که به الگوریتم این امکان را می دهد تا از تعداد تجربیات ذخیره شده مطلع شود و تصمیمات خود را براساس آن تنظیم کند.

در مجموع، این سیستم حافظه ای به الگوریتم یادگیری کمک می کند تا از تجربیات گذشته خود برای بهبود عملکرد و یادگیری بهتر استفاده کند.

توضیح عملکرد شبکه CriticNetwork

شبکه **CriticNetwork** یکی از اجزای مهم الگوریتم‌های یادگیری تقویتی است که وظیفه ارزیابی سیاست‌ها را بر عهده دارد. در این شبکه، دو تابع **Q** به طور همزمان برای تخمین ارزش وضعیت‌ها و اقدامات مختلف استفاده می‌شود. در اینجا، از دو شبکه مجزا برای هر یک از توابع **Q1** و **Q2** استفاده می‌شود تا تنوع و دقت تخمین‌ها افزایش یابد.

1. ساختار شبکه

این شبکه شامل دو بخش اصلی است:

- **Q1**: این شبکه برای تخمین مقدار **Q1** که ارزیابی اولیه از ارزش وضعیت-اقدام است طراحی شده است.
- **Q2**: شبکه دوم برای تخمین مقدار **Q2** که ارزیابی متفاوتی از همان وضعیت-اقدام را فراهم می‌کند طراحی شده است.

هر دو شبکه به طور مشابه از لایه‌های **Dense** برای پردازش اطلاعات استفاده می‌کنند:

- لایه‌های **Dense** با فعال‌سازی‌های **swish** و **relu** برای پردازش ویژگی‌های وضعیت و اقدام.
- در نهایت، خروجی نهایی به صورت تک‌بعدی برای هر شبکه **Q1** و **Q2** تولید می‌شود که نمایانگر ارزش آن وضعیت-اقدام است.

2. ورودی‌های شبکه

ورودی‌های شبکه ترکیبی از وضعیت‌ها و اقدامات هستند که به هم متصل می‌شوند:

- وضعیت (state) و اقدام (action) در طول بعد ویژگی‌ها با هم ترکیب می‌شوند تا یک ورودی مشترک برای هر دو شبکه ایجاد شود.

3. پیش‌پردازش و محاسبات

در این بخش:

- ابتدا وضعیت و اقدام با هم ترکیب می‌شوند.
- سپس از دو شبکه **Q1** و **Q2** برای پردازش این ورودی و تخمین ارزش آن استفاده می‌شود.

4. خروجی شبکه

در نهایت، این شبکه دو خروجی به شکل **q1** و **q2** تولید می‌کند که نمایانگر تخمین‌های مختلف از ارزش وضعیت-اقدام هستند. این مقادیر به الگوریتم کمک می‌کنند تا تصمیمات بهتری برای انتخاب سیاست اتخاذ کند.

توضیح عملکرد شبکه ValueNetwork

شبکه **ValueNetwork** یکی از اجزای مهم در الگوریتم‌های یادگیری تقویتی است که وظیفه تخمین تابع ارزش (Value Function) را بر عهده دارد. هدف این شبکه این است که برای هر وضعیت موجود در محیط، یک تخمین از ارزش آن وضعیت ارائه دهد که نمایانگر کیفیت آن وضعیت برای اتخاذ تصمیمات مناسب در راستای بهینه‌سازی پاداش باشد.

1. ساختار شبکه

شبکه **ValueNetwork** شامل چندین لایه **Dense** به همراه لایه‌های اضافی اختیاری برای بهبود عملکرد است:

- **fc1**: اولین لایه **Dense** با ابعاد **fc1_dims** که برای پردازش ورودی‌ها استفاده می‌شود. این لایه از تابع فعال‌سازی **ReLU** برای معرفی غیرخطی بودن به مدل استفاده می‌کند.
- **fc2**: دومین لایه **Dense** با ابعاد **fc2_dims** که پردازش‌های پیچیده‌تر را انجام می‌دهد.
- **v**: لایه نهایی که یک مقدار تک‌بعدی به عنوان خروجی شبکه تولید می‌کند. این مقدار نمایانگر تخمین ارزش وضعیت است.

2. لایه‌های اختیاری

برای بهبود پایداری فرآیند آموزش و جلوگیری از بیش‌برازش (overfitting)، دو لایه اختیاری در این شبکه وجود دارند:

- **Batch Normalization (BN):** این لایه‌ها برای نرمال‌سازی داده‌ها پس از هر لایه **Dense** استفاده می‌شوند تا پایداری فرآیند آموزش را افزایش دهند. این لایه‌ها می‌توانند به بهبود سرعت و کیفیت آموزش کمک کنند.
- **Dropout:** این لایه به‌طور تصادفی برخی از نرون‌ها را در حین آموزش غیرفعال می‌کند تا از بیش‌برازش جلوگیری کند و به مدل کمک می‌کند تا ویژگی‌های عمومی‌تری یاد بگیرد.

3. پیش‌پردازش ورودی‌ها

در این بخش، ورودی‌های شبکه که وضعیت‌های محیط هستند، به ترتیب از طریق لایه‌های مختلف پردازش می‌شوند:

- ابتدا از لایه اول **Dense** عبور می‌کنند.
- در صورت فعال بودن **Batch Normalization** و **Dropout**، این لایه‌ها به ترتیب بر روی داده‌ها اعمال می‌شوند.
- پس از پردازش در لایه دوم **Dense**، دوباره نرمال‌سازی و در صورت لزوم، **Dropout** اعمال می‌شود.
- در نهایت، خروجی نهایی به وسیله لایه **v** تولید می‌شود که یک تخمین از ارزش وضعیت است.

4. خروجی شبکه

خروجی شبکه یک مقدار تک‌بعدی است که نشان‌دهنده تخمین ارزش وضعیت ورودی است. این تخمین برای انتخاب سیاست‌های بهینه در الگوریتم‌های یادگیری تقویتی مورد استفاده قرار می‌گیرد.

توضیح عملکرد شبکه ActorNetwork

شبکه **ActorNetwork** وظیفه انتخاب سیاست در الگوریتم‌های یادگیری تقویتی را بر عهده دارد. این شبکه به‌طور خاص برای یادگیری سیاست‌های پیوسته در مسائل کنترل، مانند سیستم تعلیق خودرو، طراحی شده است. در این شبکه، هدف این است که یک توزیع نرمال از اقدامات (actions) را برای هر وضعیت (state) محاسبه کنیم، به‌طوری که مدل بتواند از آن برای انتخاب اقدام بهینه استفاده کند.

1. ساختار شبکه

شبکه **ActorNetwork** شامل چندین لایه **Dense** است که ورودی وضعیت‌ها را پردازش می‌کند:

- **fc1:** اولین لایه **Dense** با ابعاد **fc1_dims** برای پردازش ویژگی‌های وضعیت و یادگیری الگوهای غیرخطی.
- **fc2:** دومین لایه **Dense** با ابعاد **fc2_dims** برای افزایش ظرفیت شبکه و انجام پردازش‌های پیچیده‌تر.
- **mu:** لایه‌ای که میانگین توزیع نرمال (μ) را برای انتخاب اقدام پیش‌بینی می‌کند.
- **sigma:** لایه‌ای که انحراف معیار توزیع نرمال (σ) را محاسبه می‌کند، که میزان عدم قطعیت یا پراکندگی در انتخاب اقدامات را تعیین می‌کند.

2. محاسبه میانگین و انحراف معیار

شبکه ابتدا ورودی وضعیت را از طریق لایه‌های **Dense** عبور می‌دهد و سپس میانگین (μ) و انحراف معیار (σ) را به‌طور مستقل محاسبه می‌کند:

- **mu:** تخمین میانگین توزیع احتمال اقدامات.
- **sigma:** تخمین انحراف معیار توزیع، که میزان پراکندگی و تنوع در انتخاب اقدامات را تعیین می‌کند.

برای جلوگیری از ناپایداری عددی، انحراف معیار **sigma** با استفاده از تابع **exp** و مقدار **noise** (برای اطمینان از مثبت بودن آن) محاسبه می‌شود و سپس با **clip** محدود می‌شود تا از مقادیر خارج از محدوده جلوگیری شود.

3. نمونه‌برداری از توزیع نرمال (Sampling)

در این بخش، دو روش برای نمونه‌برداری از توزیع نرمال پیش‌بینی شده وجود دارد:

- **Reparameterization Trick:** این تکنیک به الگوریتم اجازه می‌دهد که از مشتقات استفاده کرده و به‌طور کارآمدتری آموزش داده شود. در این روش، با استفاده از ϵ یک متغیر تصادفی با توزیع نرمال) و انحراف معیار σ ، اقدام‌ها محاسبه می‌شوند.
- **Sampling for Inference:** در این روش، نمونه‌ها به‌طور مستقیم از توزیع نرمال با میانگین μ و انحراف معیار σ گرفته می‌شوند.

4. اعمال \tanh و محاسبه احتمال لاگ (Log Probability)

برای محدود کردن دامنه اقدامات به بازه معین، از تابع \tanh استفاده می‌شود. این اطمینان می‌دهد که اقدامات در دامنه محدود و متناسب با مقدار \max_action قرار بگیرند.

همچنین، برای محاسبه $\log probability$ اقدام‌ها، از \log احتمال توزیع نرمال استفاده می‌شود. سپس یک اصلاح کوچک با استفاده از \tanh انجام می‌شود تا احتمال لاگ درست‌تری محاسبه شود.

5. خروجی شبکه

- **Action:** اقدام انتخاب‌شده بر اساس میانگین و انحراف معیار تخمینی شبکه.
- **Log Probabilities:** احتمال لاگ اقدامات که برای محاسبات بازگشتی در الگوریتم‌های یادگیری تقویتی مورد استفاده قرار می‌گیرد (مثل به‌روزرسانی سیاست).

توضیح عملکرد کلاس Agent

کلاس **Agent** وظیفه مدیریت یادگیری و تصمیم‌گیری در یک محیط یادگیری تقویتی را بر عهده دارد. این کلاس برای پیاده‌سازی الگوریتم **Soft Actor-Critic (SAC)** طراحی شده است که از چندین

شبکه عصبی برای یادگیری سیاست، ارزیابی ارزش، و بهروزرسانی پارامترها استفاده می‌کند. در اینجا، به طور مفصل توضیح می‌دهیم که هر بخش از این کلاس چگونه عمل می‌کند.

1. مقداردهی اولیه و تنظیمات

کلاس **Agent** در ابتدا با مقادیر پیش فرض تنظیم می‌شود که شامل موارد زیر است:

- **alpha** و **beta**: این مقادیر نرخ یادگیری برای شبکه‌های **actor** و **critic** را تعیین می‌کنند.
- **gamma**: نرخ تخفیف برای پاداش‌های آینده.
- **tau**: پارامتر برای بهروزرسانی شبکه‌های هدف. (target networks)
- **memory**: حافظه برای ذخیره تجربیات، که با استفاده از کلاس **ReplayBuffer** پیاده‌سازی شده است.
- **actor, critic, value networks**: چهار شبکه عصبی برای **actor**, **critic_1**, **critic_2** و **value** تعریف می‌شوند.

2. مقداردهی شبکه‌ها

- **actor**: این شبکه برای پیش‌بینی سیاست (یعنی انتخاب اقدام‌ها) استفاده می‌شود. شبکه **actor** با توجه به وضعیت، توزیع احتمال اقدامات را تولید می‌کند.
- شبکه‌های **critic_1** و **critic_2**: این شبکه‌ها برای ارزیابی کیفیت اقدامات انتخاب‌شده توسط سیاست استفاده می‌شوند. این ارزیابی‌ها در حقیقت **Q-values** هستند که نشان‌دهنده ارزش یک اقدام خاص در یک وضعیت خاص است.
- **value**: این شبکه برای تخمین ارزش کل وضعیت‌ها استفاده می‌شود. به‌طور خاص، این شبکه تخمین می‌زند که یک وضعیت چقدر مطلوب است.
- شبکه‌های هدف (**target networks**): این شبکه‌ها نسخه‌های کپی از شبکه‌های اصلی هستند که به‌طور آهسته بهروزرسانی می‌شوند تا نوسانات در یادگیری را کاهش دهند.

3. بهروزرسانی پارامترهای شبکه‌ها

شبکه‌های **target value** و **target critic** با استفاده از پارامتر **tau** به‌طور آهسته به‌روزرسانی می‌شوند. این به‌روزرسانی‌ها با استفاده از میانگین ترکیبی از وزن‌های شبکه‌های اصلی و هدف انجام می‌شوند.

4. انتخاب اقدام

شبکه **actor** برای هر وضعیت اقدام مناسب را پیش‌بینی می‌کند. در صورت نیاز به انتخاب تصمیمات قطعی (deterministic)، فقط میانگین اقدام (μ) انتخاب می‌شود. در غیر این صورت، از **sampling** نرمال برای انتخاب اقدامات استفاده می‌شود.

5. یادگیری و به‌روزرسانی شبکه‌ها

فرآیند یادگیری به‌طور عمده شامل به‌روزرسانی شبکه‌های **critic**، **value** و **actor** است:

- **Value Network:** این شبکه برای پیش‌بینی ارزش یک وضعیت بر اساس اقدامات و-Q-values استفاده می‌شود. شبکه **value** به‌روزرسانی می‌شود تا خطای بین ارزش پیش‌بینی شده و هدف ارزشی جدید به حداقل برسد.
- **Critic Networks:** شبکه‌های **critic_1** و **critic_2** برای تخمین Q-values استفاده می‌شوند. این شبکه‌ها به‌روزرسانی می‌شوند تا خطای میانگین مربعی بین پیش‌بینی‌های فعلی و اهداف جدید حاصل از پاداش‌ها و تخمین‌های شبکه (**value**) به حداقل برسد.
- **Alpha:** پارامتر α که با استفاده از **log_alpha** نگهداری می‌شود (برای تنظیم دمای آنتروپی سیاست استفاده می‌شود. این پارامتر تأثیرگذار در میزان جستجو و تصادفی بودن اقدامات است. شبکه **alpha** به‌طور مجزا به‌روزرسانی می‌شود.

6. مدیریت تاریخچه‌ها

کلاس **Agent** اطلاعات مختلفی از روند یادگیری مانند خسارت‌های شبکه‌ها (actor)، **critic**، **value**، ضرر آلفا، و پیش‌بینی‌های مختلف شبکه‌ها را ذخیره می‌کند تا در فرآیندهای بعدی به تحلیل و بررسی عملکرد الگوریتم پرداخته شود.

7. فرآیند یادگیری

فرآیند یادگیری شامل مراحل زیر است:

- **یادگیری شبکه value:** ابتدا با استفاده از تجربیات از حافظه، هدف‌های **value** محاسبه می‌شوند و سپس شبکه **value** با استفاده از این اهداف به‌روزرسانی می‌شود.
- **یادگیری شبکه‌های critic:** سپس شبکه‌های **critic_1** و **critic_2** با استفاده از اهداف جدید به‌روزرسانی می‌شوند.
- **یادگیری شبکه alpha:** این شبکه برای تنظیم مقدار **alpha** به‌طور مستقل به‌روزرسانی می‌شود تا آنتروپی سیاست بهینه باقی بماند.

8. خروجی و ذخیره‌سازی تجربیات

تجربیات (وضعیت‌ها، اقدامات، پاداش‌ها، وضعیت‌های جدید و **done** با استفاده از متد **remember** در حافظه ذخیره می‌شوند تا در مراحل بعدی از آن‌ها برای یادگیری استفاده شود.

توضیح کلاس VehicleSuspensionEnv

کلاس **VehicleSuspensionEnv** یک محیط سفارشی برای شبیه‌سازی سیستم تعلیق خودرو در **OpenAI Gym** است. این محیط شامل ویژگی‌هایی است که به شما امکان می‌دهد سیستم تعلیق خودرو را شبیه‌سازی کرده و از آن برای یادگیری تقویتی استفاده کنید. در اینجا جزئیات هر بخش از این کلاس آورده شده است:

1. پارامترهای سیستم

در ابتدا، پارامترهای سیستم تعلیق تعریف شده‌اند. این پارامترها شامل مقادیر مختلفی هستند که رفتار سیستم تعلیق خودرو را شبیه‌سازی می‌کنند:

- **ms:** جرم قسمت‌های کشیده‌شده (sprung mass)
- **mus:** جرم بخش‌های غیرفعال (unsprung mass)
- **kf, kc, cs, cf, cc:** ضرایب سختی و میرایی مختلف سیستم.
- **kt_nl, ks_nl:** سختی‌های غیرخطی تایر و فنر.

2. فضای اقدامات (Action Space)

فضای اقدامات یک **فضای پیوسته** است که در آن یک مقدار بین -1 و 1 می‌تواند وارد شود. این مقدار با ضرب در 1000 برای اعمال نیروی کنترل به سیستم تغییر مقیاس می‌یابد.

3. فضای وضعیت (Observation Space)

فضای وضعیت 8 بعدی است که وضعیت‌های مختلف سیستم تعلیق را نمایان می‌سازد. این 8 بعد به ترتیب شامل:

- **x_us, dx_us:** موقعیت و سرعت بخش‌های غیرفعال.
- **x_s, dx_s:** موقعیت و سرعت بخش‌های کشیده‌شده.
- **x_f, dx_f:** موقعیت و سرعت فریم.
- **x_c, dx_c:** موقعیت و سرعت کوسن صندلی.

4. متد step

متد **step** برای یک گام از محیط است که در آن نیروی کنترل (بر اساس ورودی) اعمال می‌شود و وضعیت جدید سیستم محاسبه می‌شود. محاسبات مربوط به شتاب‌ها بر اساس معادلات غیرخطی صورت می‌گیرد و وضعیت جدید با استفاده از **عددگیری عددی** (numerical integration) به دست می‌آید.

- **محاسبه شتاب‌ها:** شتاب‌ها برای بخش‌های مختلف سیستم (مانند بخش‌های کشیده‌شده، بخش‌های غیرفعال و فریم) بر اساس نیروها و سختی‌ها محاسبه می‌شوند.
- **به‌روزرسانی وضعیت:** وضعیت سیستم با استفاده از شتاب‌های محاسبه‌شده به‌روزرسانی می‌شود.
- **مقیاس‌بندی وضعیت:** وضعیت به مقیاس موردنظر برای آموزش یادگیری تقویتی تبدیل می‌شود.
- **پاداش (Reward):** پاداش به‌طور منفی محاسبه می‌شود تا به مدل کمک کند تا رفتار مطلوبی پیدا کند. این پاداش بر اساس معیارهای مختلف مانند شتاب‌ها، اختلاف موقعیت‌ها، و اعمال نیروی کنترل محاسبه می‌شود.
- **پایان (Done):** در صورتی که شرایط خاصی (مانند موقعیت‌های خاص بخش‌های مختلف سیستم) برقرار شود، فرآیند به پایان می‌رسد.

5. متد reset

متد **reset** برای راه‌اندازی مجدد محیط و تنظیم وضعیت اولیه به کار می‌رود. این وضعیت اولیه معمولاً صفر است و از آن به‌عنوان نقطه شروع استفاده می‌شود.

6. متد render

متد **render** برای نمایش گرافیکی محیط استفاده می‌شود. در اینجا این متد پیاده‌سازی نشده است، اما معمولاً در شبیه‌سازی‌های پیچیده‌تر برای نمایش بصری وضعیت محیط به کار می‌رود.

7. کاربرد

این محیط برای استفاده در الگوریتم‌های یادگیری تقویتی طراحی شده است، به‌ویژه در الگوریتم‌های کنترل با بازخورد، مانند **Soft Actor-Critic (SAC)** که قبلاً در کد شما به آن اشاره شده است. با استفاده از این محیط، مدل می‌تواند سیاست بهینه‌ای برای کنترل سیستم تعلیق خودرو بیاموزد و به‌طور مؤثری به تعامل با محیط بپردازد.

آموزش شبکه

در این بخش، یک حلقه آموزشی برای یک عامل یادگیری تقویتی طراحی شده است که در محیط شبیه‌سازی سیستم تعلیق خودرو فعالیت می‌کند. هدف این حلقه، آموزش عامل برای پیدا کردن سیاست بهینه‌ای است که به آن اجازه می‌دهد تا نیروی کنترل مناسبی را به سیستم اعمال کرده و پاداش بیشتری دریافت کند. مراحل آموزشی به‌صورت زیر انجام می‌شود:

1. شروع هر قسمت: (Episode)

- در ابتدای هر قسمت، وضعیت محیط (state) با استفاده از متد **reset** بازیابی می‌شود.
- وضعیت اولیه به عامل داده می‌شود و قسمت شروع می‌شود.

2. عملیات داخل قسمت: (Episode)

- در هر گام از قسمت، عامل یک اقدام (action) انتخاب می‌کند. این اقدام بر اساس سیاستی است که عامل در طول زمان یاد می‌گیرد.
- عامل به محیط ارسال می‌کند که چه اقدامی انجام دهد و سپس محیط وضعیت جدید، پاداش و وضعیت اتمام را باز می‌گرداند.
- عامل این اطلاعات را ذخیره می‌کند تا در آینده از آن برای به‌روزرسانی شبکه‌های عصبی خود استفاده کند.

3. یادگیری عامل:

- پس از انجام هر گام، عامل از داده‌های ذخیره‌شده استفاده می‌کند تا شبکه‌های عصبی خود (که نماینده سیاست و ارزش‌ها هستند) را به‌روزرسانی کند.
- یادگیری شامل محاسبه و به‌روزرسانی ضررهای مربوط به شبکه‌های مختلف (مقدار، منتقد و آلفا) است.

4. محاسبه پاداش:

- پاداش در هر گام از قسمت محاسبه می‌شود و به عامل داده می‌شود. این پاداش از ترکیب فاکتورهای مختلف مانند انحرافات از مقادیر هدف، شتاب‌ها و میزان نیروی کنترل است.

5. پیگیری تاریخچه عملکرد:

- برای ارزیابی پیشرفت عامل، نمره هر قسمت ذخیره می‌شود.

- به طور دوره‌ای (هر 5 قسمت)، میانگین نمرات 100 قسمت آخر محاسبه می‌شود تا عملکرد عامل در طول زمان ارزیابی شود.

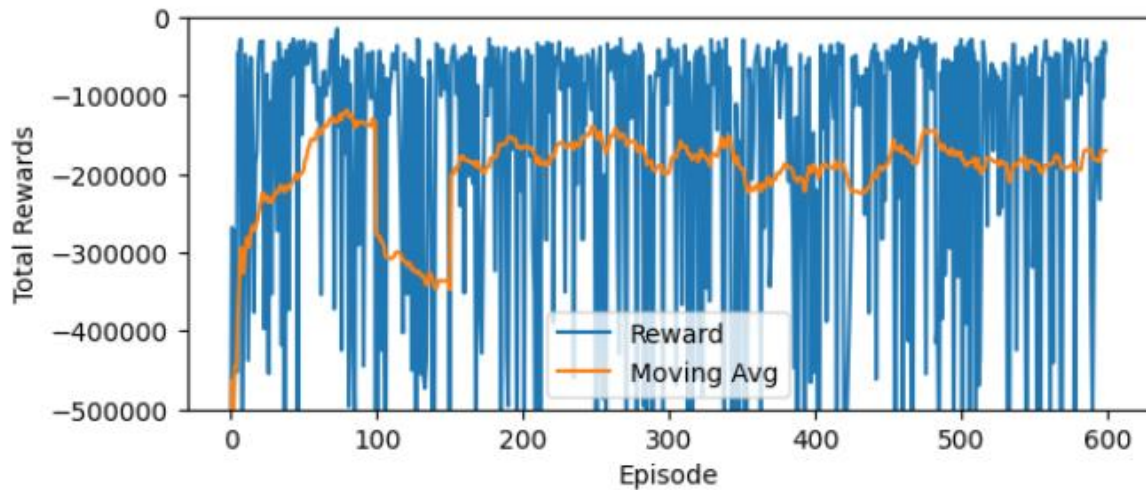
6. پایان قسمت:

- قسمت زمانی به پایان می‌رسد که شرایط خاصی مانند اتمام زمان یا رسیدن به وضعیت غیرمطلوب (مثل انحراف‌های بیش از حد از مقادیر هدف) محقق شود.

هدف آموزش:

هدف این فرایند، بهینه‌سازی سیاست عامل است تا با کاهش پاداش منفی (که ناشی از شتاب‌های بالا، انحرافات از هدف، و نیروی کنترل زیاد است) سیستم تعلیق خودرو به طور مؤثری کنترل شود. این مدل به مرور زمان یاد می‌گیرد که بهترین اقدامات را در شرایط مختلف برای کاهش پاداش منفی و بهبود عملکرد سیستم انتخاب کند.

نتایج آموزش



شکل 5-1 بررسی پاداش و حرکت در طی آموزش

همانطور که در نمودار فوق مشاهده میکنید در طول آموزش میزان پاداش ها بشدت نوسانی میباشد اما حرکت به میزان نسبتا مناسبی در برخی اپیزودها پیشرفت داشته است

دلیل نوسانات شدید در میزان پاداش ها و پیشرفت نسبی در برخی اپیزودها می تواند به عوامل مختلفی برگردد که به روند یادگیری عامل و پیچیدگی محیط مرتبط هستند. این نوسانات و پیشرفت ها ممکن است ناشی از موارد زیر باشند:

1. آغاز آموزش و کاوش در محیط:

- در مراحل اولیه آموزش، عامل بیشتر در حال کاوش در محیط است تا یادگیری یک سیاست بهینه. در این دوره، عامل ممکن است تصمیمات اشتباهی بگیرد که منجر به پاداش های منفی و نوسانات بالا در پاداش ها شود.

- در این مرحله، عامل هنوز استراتژی مناسبی برای مواجهه با محیط پیدا نکرده است و در نتیجه ممکن است در برخی اپیزودها پاداش‌های خیلی پایین یا حتی منفی دریافت کند. این امر باعث نوسانات زیاد در نمودار پاداش‌ها می‌شود.

2. آموزش با سیاست‌های تصادفی:

- عامل در ابتدا با سیاست تصادفی اقدام می‌کند، که منجر به پاداش‌های متغیر و نوسانات شدید در هر اپیزود می‌شود. برخی اپیزودها ممکن است به‌طور تصادفی بهتر از دیگران عمل کنند و منجر به پاداش‌های بالاتر شوند، در حالی که در برخی دیگر عملکرد ضعیفی از خود نشان می‌دهد.
- این پدیده باعث ایجاد نوسانات در پاداش‌ها می‌شود و در عین حال، عوامل تصادفی می‌توانند باعث ظهور نتایج تصادفی و غیرقابل پیش‌بینی شوند.

3. الگوریتم یادگیری و زمان‌بندی به‌روزرسانی‌ها:

- الگوریتم‌های یادگیری تقویتی معمولاً از به‌روزرسانی‌های تدریجی استفاده می‌کنند. این به‌روزرسانی‌ها در ابتدا ممکن است خیلی کوچک و ناپایدار باشند و نوسانات زیادی در پاداش‌ها ایجاد کنند.
- یادگیری در ابتدا ممکن است در دنیای واقعی به کندی پیش برود، زیرا عامل هنوز بهینه‌ترین سیاست را پیدا نکرده است. تغییرات کوچک در پارامترهای شبکه عصبی و اشتباهات جزئی در انتخاب‌های اولیه می‌توانند باعث نوسانات در پاداش‌ها شوند.

4. چالش‌های محیطی و پیچیدگی دینامیک‌های سیستم:

- محیط پیچیده‌ای مانند سیستم تعلیق خودرو می‌تواند موجب ایجاد نوسانات شود. به دلیل وجود پارامترهای غیرخطی و پیچیدگی‌های فیزیکی مانند شتاب‌ها، نیروهای کنترل و تعاملات پیچیده میان اجزاء مختلف سیستم، عامل ممکن است نتواند در همه شرایط عملکرد یکنواختی داشته باشد.

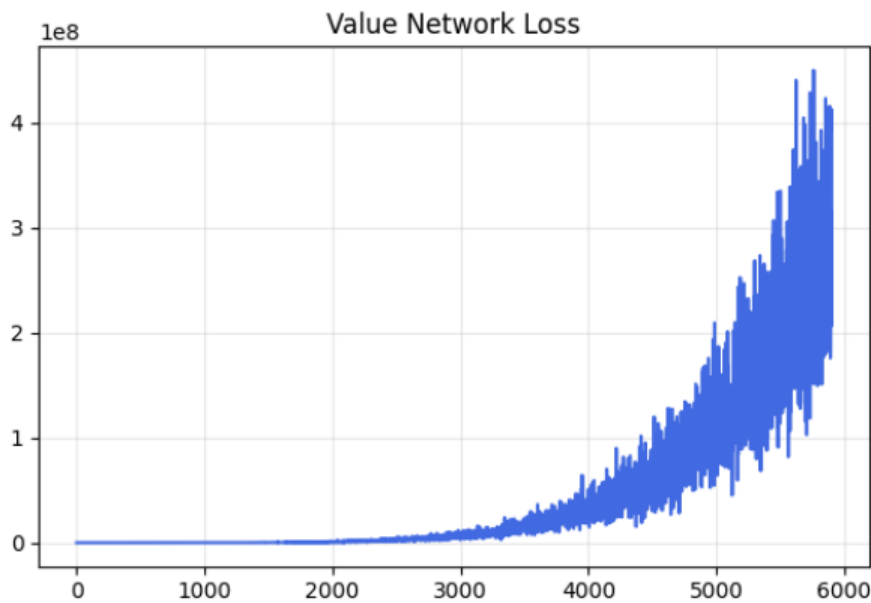
- در برخی اپیزودها، به‌ویژه زمانی که سیستم نیاز به واکنش سریع به تغییرات دارد، عامل ممکن است به‌طور غیرمنتظره‌ای به نتایج مثبت دست یابد و در برخی دیگر، به دلیل محدودیت‌های مدل و رفتار غیرقابل پیش‌بینی سیستم، عملکرد ضعیف‌تری داشته باشد.

5. نوسانات در پاداش به دلیل پاداش‌دهی پیچیده:

- پاداش‌دهی که بر اساس معیارهایی مانند انحرافات از هدف یا شتاب‌های بالا صورت می‌گیرد، ممکن است در ابتدا باعث نوسانات زیاد شود. به‌ویژه اگر عامل هنوز سیاست مناسبی برای تعادل بین این معیارها پیدا نکرده باشد.
- این نوسانات می‌توانند در اثر عدم تطابق کامل با اهداف پاداش‌دهی در آغاز آموزش رخ دهند. برای مثال، ممکن است عامل در برخی اپیزودها در تلاش برای کاهش شتاب‌های بزرگ، به نتایج بهتری دست یابد، در حالی که در اپیزودهای دیگر به دلیل تمرکز بیش از حد بر روی جنبه‌های خاص سیستم، از عملکرد بهینه دور شود.

6. پیشرفت‌های تدریجی در یادگیری:

- با گذر زمان، عامل شروع به یادگیری و بهبود سیاست خود می‌کند، اما این پیشرفت‌ها به‌طور یکنواخت اتفاق نمی‌افتد. در برخی اپیزودها، پیشرفت قابل توجهی در یادگیری سیاست ایجاد می‌شود که موجب بهبود پاداش‌ها می‌شود.
- در این مراحل، عامل ممکن است درک بهتری از محیط پیدا کرده و عملکردش را به‌طور چشمگیری بهبود دهد، که باعث می‌شود در برخی از اپیزودها پیشرفت قابل توجهی مشاهده شود.

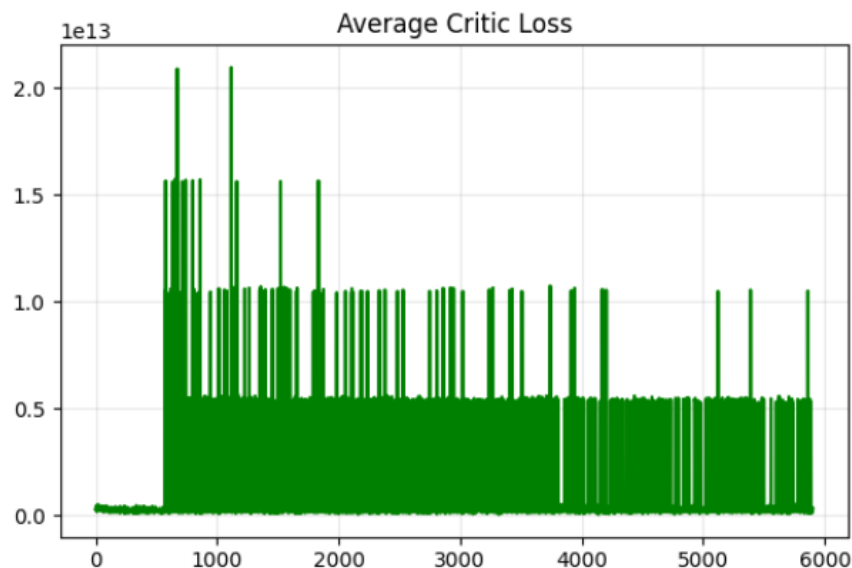


شکل 5-2 میزان هزینه شبکه

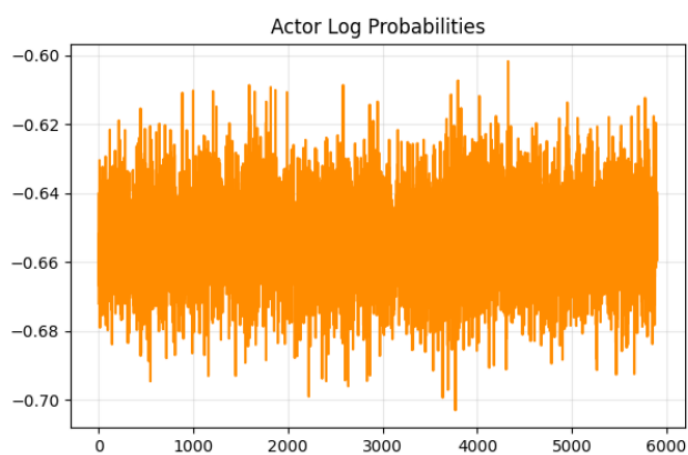
نمودار هزینه (Value Loss) نشان می‌دهد که میزان خطای پیش‌بینی ارزش وضعیت‌ها به‌طور نامنظم و ناپایدار افزایش می‌یابد. این افزایش ناپایداری ممکن است به دلایل متعددی رخ دهد:

- **عدم ثبات در یادگیری:** در مراحل اولیه یا در نقاطی که عامل هنوز به سیاست بهینه دست نیافته، تغییرات ناگهانی در وزن‌های شبکه ممکن است باعث افزایش موقت هزینه شود. این موضوع به‌ویژه در سیستم‌های پیچیده با دینامیک‌های غیرخطی مشهود است.
- **تنوع داده‌های ورودی:** تغییر وضعیت‌ها و تجربیات متفاوت در محیط، مدل ارزش ممکن است با توزیع‌های متنوعی از داده مواجه شود که در نتیجه به‌روزرسانی‌های شبکه به‌طور نامنظم انجام می‌شود و هزینه به صورت ناپایدار تغییر می‌کند.

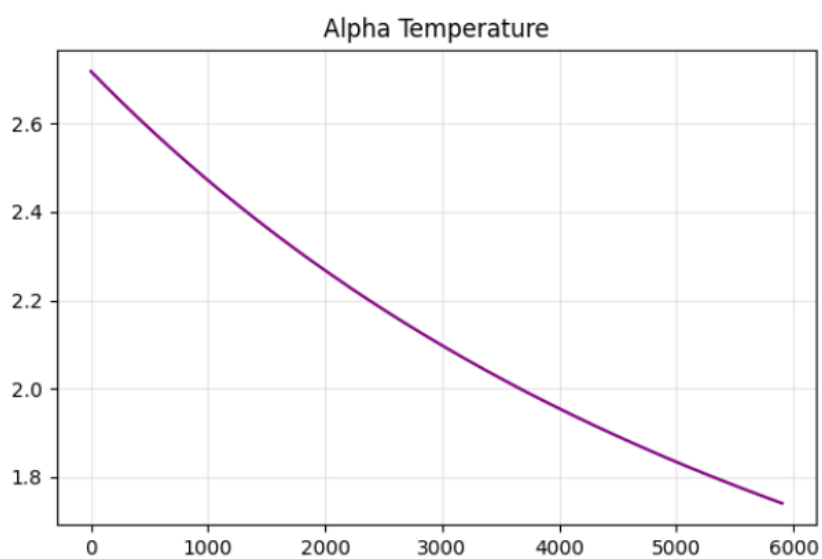
- **تنظیمات نرخ یادگیری و پارامترها:** استفاده از نرخ یادگیری نسبتاً بالا یا به‌روزرسانی‌های ناگهانی پارامترها می‌تواند باعث شود که مدل به بهبود تدریجی نرسد و هزینه به طور ناگهانی افزایش پیدا کند.
 - **اثر اکتشاف (Exploration):** در فرآیند یادگیری، عامل ممکن است به دلیل استراتژی‌های کاوش، گاهی اقداماتی انجام دهد که باعث می‌شود مقدار پیش‌بینی شده ارزش با هدف واقعی فاصله بگیرد. این موضوع می‌تواند به صورت موقت باعث افزایش هزینه شود.
- در مجموع، این نوسانات در نمودار هزینه بیانگر چالش‌های موجود در تخمین صحیح ارزش وضعیت‌ها و بهبود تدریجی مدل در مواجهه با محیط‌های پیچیده هستند. به مرور زمان و با بهبود سیاست عامل، انتظار می‌رود که این نوسانات کاهش یافته و هزینه به یک سطح پایدارتر نزدیک شود.



شکل 3-5 میانگین critic loss



شکل 4-5 میزان احتمالات اکتور



شکل 5-5 میزان تغییر alpha temperature

در الگوریتم SAC، پارامتر دمای α نقش کلیدی در تنظیم تعادل بین اکتشاف (exploration) و بهره‌برداری (exploitation) ایفا می‌کند. کاهش یکنواخت این دما در نمودار نشان‌دهنده روند تطبیق خودکار این پارامتر بر اساس سیاست یادگرفته‌شده است. به عبارت دیگر:

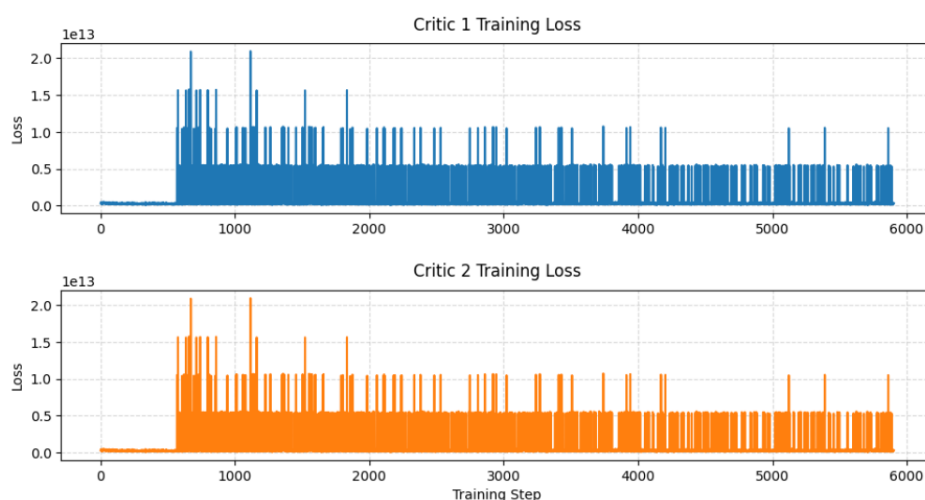
تنظیم خودکار آنتروپی: در ابتدای آموزش، به دلیل نیاز به اکتشاف گسترده، دمای α ممکن است بالاتر باشد تا عامل بتواند با سیاست تصادفی، محیط را به خوبی کاوش کند.

کاهش تدریجی با بهبود سیاست: با پیشرفت یادگیری، عامل به تدریج سیاست خود را به سمت تصمیم‌گیری‌های هدفمندتر و قطعی‌تر هدایت می‌کند. در این حالت، مقدار آنتروپی (عدم قطعیت) کاهش می‌یابد و نیاز به میزان بالای اکتشاف کمتر می‌شود.

تنظیم دما به سمت هدف آنتروپی: در فرآیند به‌روزرسانی α ، اگر مقدار آنتروپی سیاست از مقدار هدف (target entropy) بالاتر باشد، به‌روزرسانی گرادینان منجر به کاهش دمای α می‌شود تا سطح آنتروپی به مقدار مطلوب نزدیک شود.

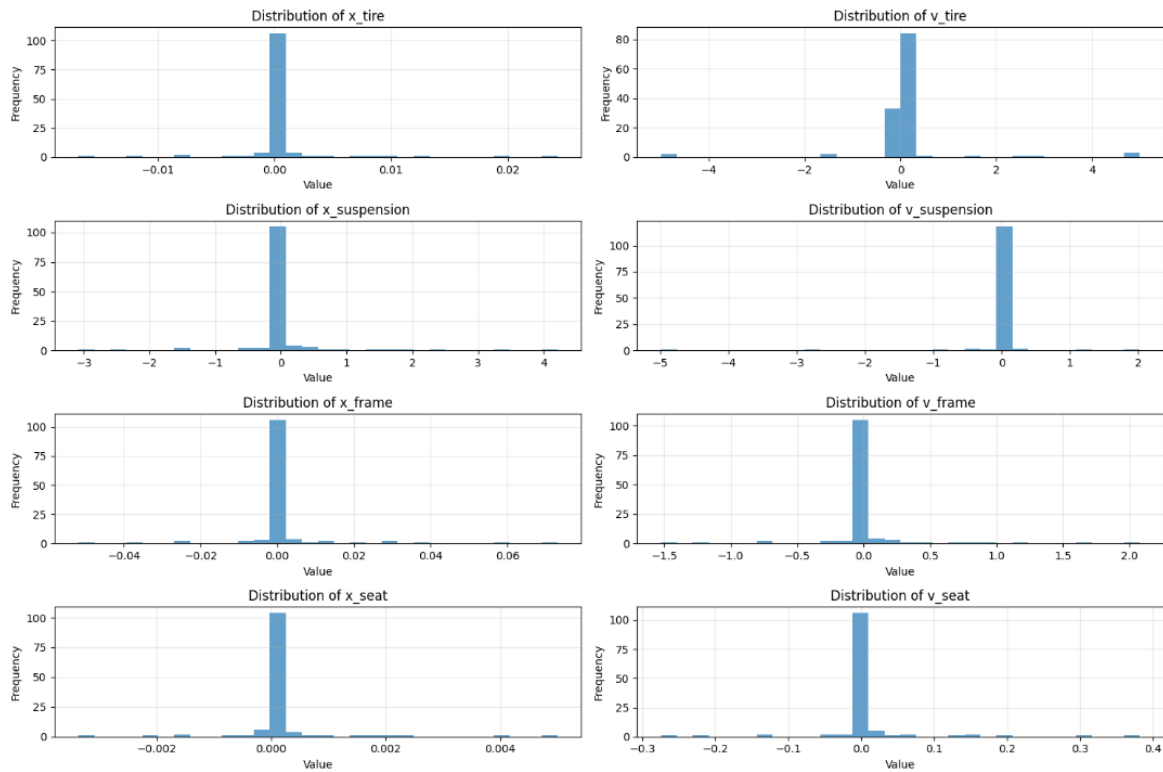
بنابراین، کاهش یکنواخت دمای α نشان‌دهنده‌ی این است که عامل به مرور زمان به سیاستی نزدیک می‌شود که سطح آنتروپی آن مطابق با هدف تعیین‌شده است و نیازی به اکتشاف گسترده کمتری دارد. این روند، نشانه‌ای از پایداری و تثبیت سیاست یادگرفته‌شده در طول آموزش محسوب می‌شود.

SAC Critic Networks Training Progress



شکل 5-6 تغییرات critic loss

در نمودار critic loss می‌بینیم که هزینه شبکه‌های منتقد نسبت به ابتدای آموزش کاهش یافته است. این بهبود نشان‌دهنده افزایش دقت تخمین‌های Q-value است؛ یعنی شبکه‌های منتقد توانسته‌اند به مرور زمان با استفاده از تجربیات جمع‌آوری‌شده، ارزش وضعیت-اقدام را بهتر پیش‌بینی کنند. کاهش critic loss بیانگر تثبیت و بهبود عملکرد سیستم ارزیابی و در نتیجه پیشرفت کلی در بهینه‌سازی سیاست عامل است.



شکل 5-7 مانیتور فرکانس هر حالت

این نمودار نشان می‌دهد که مقادیر وضعیت‌ها عمدتاً حول نقطه تعادل (equilibrium) یا فرکانس تعادل متمرکز هستند. به عبارت دیگر، سیستم در طول زمان به حالت پایدار نزدیک شده و نوسانات زیاد کاهش یافته‌اند. این موضوع بیانگر کارایی بهبود یافته سیستم کنترل است، به گونه‌ای که عامل قادر است رفتار سیستم را به نحوی تنظیم کند که وضعیت‌ها حول نقطه تعادل باقی بمانند و از تغییرات ناخواسته و شدید جلوگیری شود.

نتیجه گیری کلی

در این آزمایش، با بهره‌گیری از الگوریتم یادگیری تقویتی SAC، عامل قادر شد تا کنترل نسبتاً خوبی بر سیستم تعلیق خودرو اعمال کند. نتایج به دست آمده نشان می‌دهد که به مرور زمان، وضعیت‌های سیستم حول نقطه تعادل متمرکز شده و نوسانات زیاد کاهش یافته‌اند. کاهش هزینه شبکه‌های منتقد و شبکه ارزش، بیانگر بهبود دقت تخمین عملکرد سیستم و تثبیت سیاست یادگرفته‌شده عامل است.

از یک سو، این نتایج نشان می‌دهد که عامل توانسته است با استفاده از تکنیک‌های پیشرفته یادگیری تقویتی، به درک بهتری از دینامیک سیستم دست یابد و کنترل مناسبی را ارائه دهد. از سوی دیگر، باید توجه کرد که سیستم تعلیق خودرو دارای پیچیدگی‌های بسیار بالایی است؛ به دلیل وجود دینامیک‌های غیرخطی، اثرات نیروهای مختلف و تعامل پیچیده اجزاء مختلف سیستم، دستیابی به یک کنترل کاملاً بهینه چالشی جدی محسوب می‌شود.

برای بهبود عملکرد کنترل و رسیدن به نتایج مطلوب‌تر، می‌توان از راهکارهای زیر استفاده کرد:

- **بهینه‌سازی ساختار شبکه‌های عصبی:** استفاده از معماری‌های عمیق‌تر یا اعمال تکنیک‌های پیشرفته تنظیم وزن می‌تواند به بهبود دقت تخمین‌ها و کاهش خطاهای یادگیری کمک کند.
- **تنظیم دقیق هایپرپارامترها:** بهینه‌سازی نرخ‌های یادگیری، پارامترهای آلفا (دمای آنتروپی) و پارامترهای مرتبط با به‌روزرسانی شبکه‌های هدف می‌تواند باعث پایداری بیشتر در فرایند آموزش شود.
- **استفاده از تکنیک‌های پیشرفته کنترل:** ترکیب الگوریتم‌های یادگیری تقویتی با روش‌های کنترل تطبیقی یا الگوریتم‌های تکاملی می‌تواند به دستیابی به سیاست‌های کنترل بهینه‌تر و تطبیق بهتر با تغییرات دینامیکی سیستم کمک کند.
- **افزایش کیفیت داده‌های آموزشی:** جمع‌آوری تجربیات بیشتر و استفاده از استراتژی‌های پیشرفته بازپخش تجربیات می‌تواند به کاهش نوسانات و افزایش پایداری سیستم منجر شود.

در مجموع، اگرچه کنترل نسبتاً مناسبی در این آزمایش به دست آمده است، اما با توجه به پیچیدگی بیش از حد سیستم، همچنان فضای بزرگی برای بهبود و پیشرفت وجود دارد. با اعمال راهکارهای پیشنهادی و

بهینه‌سازی بیشتر مدل، می‌توان انتظار داشت که عملکرد عامل در کنترل سیستم تعلیق به صورت چشمگیری ارتقا یابد.

فصل ششم:

جمع‌بندی و نتیجه‌گیری

نتیجه‌گیری روش‌های مختلف

در این پروژه، روش‌های هوشمند مختلف برای کنترل و عیب‌یابی سیستم تعلیق خودروی پیشرفته با استفاده از مدل‌های مختلف شبکه عصبی و الگوریتم‌های یادگیری تقویتی مورد بررسی قرار گرفتند. نتایج حاصل از هر روش به‌طور جداگانه تحلیل شده و در نهایت مقایسه‌ای دقیق بین این روش‌ها انجام خواهد گرفت.

MLP (Perceptron چندلایه)

شبکه عصبی چندلایه (MLP) به‌عنوان یک مدل ساده و قدرتمند برای شبیه‌سازی سیستم‌های پیچیده در این پروژه مورد استفاده قرار گرفت. نتایج نشان‌دهنده عملکرد بسیار خوب این مدل در شناسایی و تعقیب دینامیک سیستم تعلیق خودروی پیشرفته بود. MLP توانست به‌خوبی روابط پیچیده در داده‌ها را شبیه‌سازی کرده و دقت قابل قبولی در پیش‌بینی متغیرهای حالت سیستم ارائه دهد. همچنین، به‌دلیل ساختار ساده‌تر و نیاز به منابع محاسباتی کمتر، این روش به‌عنوان گزینه‌ای مناسب برای پیاده‌سازی سریع و کارآمد در سیستم‌های مشابه معرفی می‌شود.

RBF (شبکه عصبی شعاعی)

در استفاده از شبکه عصبی RBF برای شناسایی سیستم، نتایج حاکی از مشکلاتی در عملکرد این مدل بود. به‌ویژه، پس از اعمال مدل بر روی داده‌های سیستم تعلیق، مشاهده شد که خروجی‌های مدل در مقایسه با پیش‌بینی‌های مورد انتظار همچنان ناپایدار و غیرمطلوب بودند. این ناپایداری در عمل باعث کاهش دقت کنترل سیستم و حتی ناتوانی در شبیه‌سازی دینامیک واقعی سیستم شد. از آنجایی که سیستم تعلیق خودرو از دینامیک‌های پیچیده و غیرخطی برخوردار است، این روش قادر به شبیه‌سازی این پیچیدگی‌ها به‌درستی نبود و به‌عنوان یک گزینه ناکارآمد در نظر گرفته شد.

LSTM (شبکه‌های عصبی بازگشتی طولانی‌مدت کوتاه‌مدت)

در این بخش، از شبکه‌های عصبی LSTM برای شناسایی و ریت سیستم تعلیق استفاده شد. نتایج این مدل نشان داد که LSTM نمی‌تواند به‌طور مناسب دینامیک سیستم تعلیق خودروی پیشرفته را مدل‌سازی کند. به‌ویژه در فرآیند آموزش، میزان خطای مدل بسیار بالا بود و مدل نتوانست همگرایی لازم

را به‌دست آورد. این مشکلات ناشی از ضعف در تعمیم‌پذیری مدل و ناکارآمدی آن در تخمین متغیرهای حالت سیستم بود. به‌دلیل این نتایج ضعیف، پیشنهاد می‌شود که از روش‌های جایگزین مانند فیلتر کالمن توسعه‌یافته (EKF) یا ترکیب شبکه‌های عصبی با مدل‌های فیزیکی (Hybrid Models) استفاده شود تا دقت و پایداری بالاتری در شبیه‌سازی دینامیک سیستم حاصل گردد.

SAC (الگوریتم یادگیری تقویتی Soft Actor-Critic)

در این بخش، از الگوریتم یادگیری تقویتی SAC برای کنترل سیستم تعلیق خودرو استفاده شد. نتایج به‌دست‌آمده نشان داد که این الگوریتم توانسته است به‌طور موفقیت‌آمیزی کنترل مناسبی بر سیستم اعمال کند. عامل یادگیری تقویتی SAC در طی فرآیند آموزش توانست نوسانات سیستم را کاهش داده و به‌طور تدریجی وضعیت‌های سیستم را به نقطه تعادل نزدیک کند. با این حال، به‌دلیل پیچیدگی‌های بالای سیستم تعلیق با چهار درجه آزادی، همچنان چالش‌هایی در دستیابی به کنترل بهینه و پایدار وجود دارد. به‌طور کلی، این روش می‌تواند عملکرد بهتری داشته باشد اگر که با استفاده از تکنیک‌های پیشرفته‌تر مانند تنظیم دقیق‌های پیرامون‌ها یا ترکیب الگوریتم‌های یادگیری تقویتی با روش‌های کنترل تطبیقی بهبود یابد.

مقایسه روش‌ها

معیار	MLP	RBF	LSTM	SAC
دقت	بالا	پایین	ضعیف	متوسط
پایداری سیستم کنترل	پایدار	ناپایدار	ناپایدار	نسبتاً پایدار
توانایی در شبیه‌سازی پیچیدگی‌های غیرخطی	خوب	ضعیف	ضعیف	متوسط
زمان آموزش	کوتاه	متوسط	طولانی	طولانی
پایداری در طول زمان	پایدار	ناپایدار	ناپایدار	پایدار
نیاز به منابع محاسباتی	کم	متوسط	زیاد	زیاد
کاربرد در سیستم‌های پیچیده	مناسب	نامناسب	نامناسب	مناسب

جدول 6-1 مقایسه روش‌های مختلف

نتیجه‌گیری کلی

با توجه به نتایج به‌دست‌آمده از بررسی روش‌های مختلف، می‌توان نتیجه گرفت که برای سیستم‌های پیچیده‌ای مانند سیستم تعلیق خودرو با دینامیک‌های غیرخطی، هر یک از این روش‌ها دارای مزایا و معایب خاص خود هستند. روش MLP به‌عنوان بهترین گزینه در شبیه‌سازی و شناسایی دینامیک سیستم‌های غیرخطی پیچیده با دقت بالا و پایداری مناسب شناخته شد. الگوریتم SAC نیز در کنترل سیستم تعلیق عملکرد مطلوبی داشت، اما نیاز به بهینه‌سازی و تنظیم دقیق‌تر دارد. در مقابل، مدل‌های RBF و LSTM در شبیه‌سازی و کنترل سیستم با مشکلاتی همراه بودند که کارایی آن‌ها را به‌شدت کاهش داد.

در نهایت، پیشنهاد می‌شود برای سیستم‌های مشابه از ترکیب روش‌های مختلف یادگیری عمیق و یادگیری تقویتی استفاده شود. استفاده از معماری‌های پیشرفته‌تر، بهینه‌سازی‌های بیشتر و ترکیب روش‌های کنترل کلاسیک با یادگیری عمیق می‌تواند به دستیابی به عملکرد بهینه‌تر و پایدارتر منجر شود.

منابع و مراجع

- Optimization of nonlinear quarter car suspension– seat–driver model [1]
Mahesh P. Nagarkar a,b, Gahininath J. Vikhe Patil c , Rahul N. Zaware Patil
- Simulation and Analysis of Passive and Active Suspension System Using [2]
Quarter Car Model for Different Road Profile Abdolvahab Agharkakli ,
Ghobad Shafiei Sabet , Armin Barouz