

Megerősítéses tanulás alapú robusztus járműirányítás autonóm járművek pályakövetésére

Attila Lelkó*. Balázs Németh*.

*HUN-REN Számítástechnikai és Automatizálási Kutatóintézet
1111, Budapest Kende u. 13-17.

Absztrakt: Ezen publikáció egy új módszert mutat be, mellyel megerősítéses tanulás alapú irányítási módszerek ötvözhetőek klasszikus robusztus irányítási metódusokkal. A kombináció révén egy magas minőségi kritériumokat teljesítő robusztus irányítási rendszer adódik. Az ismertetett módszer egy autonóm jármű irányításán keresztül kerül bemutatásra. A megerősítéses tanulás során választott jutalom függvény megválasztásával különféle vezetési stílusok valósíthatók meg, pl. köríró minimalizálás, pályakövetés, utazási kényelem. A neurális hálózat tanítása a Proximal Policy Optimization algoritmussal történt, a robusztus irányítás pedig \mathcal{H}_∞ alapú. A két szabályzó egy felügyelő struktúra segítségével kerül kombinálásra, melyben egy kvadratikusan optimalizálási feladat valósul meg. A módszer eredményeként egy olyan irányítási struktúra adódik, mely a jármű hossz- és oldalirányú irányítását is megvalósítja a referenciasebesség és a kormánysszög előírásával. Az algoritmus hatékonysága szimulációkon keresztül kerül bemutatásra.

1. BEVEZETÉS

Manapság, a számítástechnika fejlődése nyomán megjelenő gyorsabb és gyorsabb számítógépes eszközök, illetve hardverek megjelenésével a gépi tanulás alapú módszerek is gyors fejlődésnek indultak és jelentek meg szinte minden tudományterületen, többek között a komplex feladatokat megoldó irányítórendszerekben is. Egy tipikus példa erre az autonóm járművek irányítása, ahol különböző környezetérzékelési, döntéshozási és irányítási problémákat kell megoldani egy folyamatosan változó közlekedési környezetben. Ilyen problémák esetében a magas performanciájú és megbízható irányítási megoldások tervezése komoly kihívás rejt.

Egy lehetséges megoldást jelenthetnek a klasszikus irányítási eljárások adattokkal támogatott kiterjesztései, például Model Predictive Control (MPC, Kabzan et al. (2019); McKinnon és Schoellig (2019); Rosolia és Borrelli (2020)), model-free (modell mentes) irányítás (MFC, Fliess és Join (2021); Fényes et al. (2022)), robusztus és lineáris paraméter változó módszerek (Németh and Gáspár (2021); Sename (2021)). Ezen felül robusztus irányítási módszerek is relevánsak az autonóm járművek kontextusában, hiszen ezeknek a rendszerek általában különféle zajokkal, zavarásokkal és bizonytalanságokkal terhelték. Szerencsére azonban ezek kezelhetőek klasszikus módszerekkel, mint például Khosravani et al. (2014) publikációjában, ahol a robusztus irányítás egy driver-in-the-loop esetre lett tervezve.

Egy másik lehetséges megközelítése a magas minőségi kritériumokat igénylő irányítási feladatoknak a gépi tanulás alapú eljárások alkalmazása. Ezek közül a legerterjedtebbek általában valamilyen mély neurális hálózaton alapulnak. Ezen

algoritmusok előnye, hogy jelentősen képesek profitálni a valós viselkedés során gyűjtött nagy mennyiségű adatból, ezáltal megvalósítva az optimális viselkedést (Brüggenmann and Possieri (2021)). Különféle megoldások találhatók az irányítástechnikában, ahol a szabályzót egy neurális hálózat valósítja meg, ezeknek hatékonyságára lehet példa He et al. (2018), Haiyang et al. (2016) és Zhang et al. (2021). A hátránya ezeknek a módszereknek önmagában azonban, hogy hiányoznak az elméleti garanciák a viselkedésükre általános esetben, ez különösen fontos lenne biztonságkritikus alkalmazásokban, mint amilyenek a járműirányítási rendszerek is. A neurális hálózat alapú rendszerek validációjára tesz kísérletet több publikáció is, mint például Lelkó et al. (2021), ahol a rendszer különböző munkapontokban történő linearizálásával egy politopikus rendszer képződik, mely már Lyapunov módszerrel vizsgálható. További lehetséges megközelítéseket ismertetnek Ruan et al. (2018), Huang et al. (2017) és Wu et al. (2020).

Ezen publikáció célja, hogy bemutasson egy integrált járműirányítási stratégiát, amellyel a klasszikus model alapú irányítási módszerek és a modern gépi tanulás alapú módszerek előnyei kombinálhatók. A kombinált megoldás képes a gépi tanulási módszerekhez hasonló magas performanciára, miközben rendelkezik a klasszikus módszerek megbízhatóságával és robusztusságával.

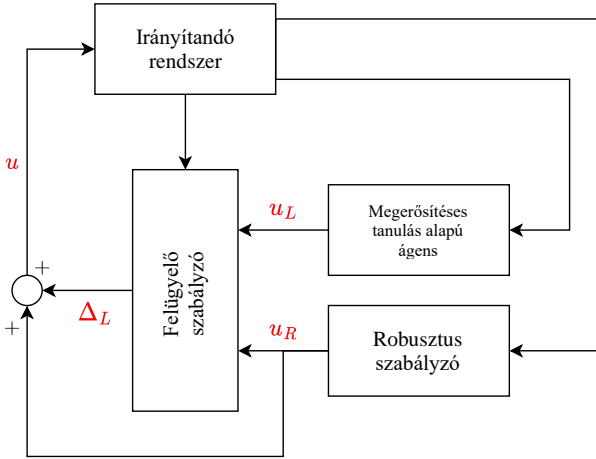
2. KOMBINÁLT IRÁNYÍTÁSI STRUKTÚRA

A klasszikus algoritmusok és a gépi tanulás alapú módszerek integrálása egy felügyelő szabályzó segítségével történt. Az így adódó struktúra látható az 1. ábrán. Az ismertetett struktúrában a robusztus szabályzó és a megerősítéses tanulás alapú ágens (RL ágens) egymástól függetlenül működik és

számolja ki a lehetséges beavatkozó jelet minden időpillanatban. A tényleges beavatkozó jelet azonban a felügyelő szabályzó határozza meg, figyelembe véve az előbbi két szabályzó kimenetét. A beavatkozó jel meghatározása egy additív ΔL érték segítségével történik a következőképpen:

$$u = u_R + \Delta L, \quad (1)$$

ahol u_R a robusztus irányítás kimenetét jelöli. ΔL egy korlátos additív jel, melyet a felügyelő szabályzó számít ki, értékei ΔL_{min} és ΔL_{max} között változhatnak.



1. ábra. A felügyelő szabályzó struktúra blokkdiagramja.

A felügyelő szabályzó célja a neurális hálózat által javasolt és a ténylegesen megvalósuló beavatkozás közötti különbség minimalizálása, figyelembe véve a biztonságos mozgáshoz szükséges korlátozásokat. A robusztus stabilitást a robusztus irányítás biztosítja. Klasszikus módszerek, mint például a \mathcal{H}_∞ alapú irányítástervezés képesek kezelni a korlátos bemeneti zavarásokat, képesek garantálni a robusztus stabilitást és performanciát korlátos bemeneti zavarás esetén is.

A felügyelő szabályzó által megoldott optimalizálási feladat a következőképpen néz ki:

$$\operatorname{argmin}_{\Delta L} \| u_R + \Delta L - u_L \|_2^2 \quad (2a)$$

feltéve, hogy

$$\Delta L \in [\Delta L_{min}; \Delta L_{max}] \quad (2b)$$

$$e_{Lat,T} \leq e_{max}, \quad (2c)$$

ahol $e_{Lat,T}$ a jármű modellje alapján becsült pályakövetési hiba. A becslést a felügyelő szabályzóba beépített modell alapú predikciós réteg végzi. Ezzel az optimalizálási feladattal nem csak a robusztus performanciához szükséges bemeneti additív zavarás korlátossága teljesül, hanem egyéb követelmények is figyelembe vehetők, ebben az esetben például a pályakövetési hiba korlátozása, mellyel megakadályozható a pálya elhagyása.

3. SZABÁLYZÓK TERVEZÉSI FOLYAMATA

Ebben a fejezetben az ismertetett struktúrában szereplő két szabályzó, a robusztus irányítás és a megerősítéses szabályzó alapú ágens tervezésének folyamata kerül ismertetésre. Mindkettő irányítás kritikus részét képezi a teljes struktúrának, hiszen a felügyelő szabályzó működése ezeken alapszik. A robusztusságért a robusztus szabályzó, míg a magas performanciáért a megerősítéses tanulás alapú ágens felelős.

3.1 Robusztus irányítás tervezése

Ebben az alfejezetben a robusztus \mathcal{H}_∞ irányítástervezés kerül áttekintésre a jármű laterális irányítása esetén, illetve a biztonságos referenciasebesség megválasztása kerül ismertetésre.

Az irányítástervezés a rendszer egyszerűsített modelljén alapszik, ami lehet például egy egyszerű dinamikus biciklimodell Kong et al. (2015) alapján. A linearizált dinamikus biciklimodell egyszerűen állapotteres alakra hozható, és felhasználva (1)-es összefüggést:

$$\dot{x} = Ax + Bu = Ax + Bu_R + B\Delta L, \quad (3)$$

ahol A és B a rendszer mátrixai, $x = [\dot{y} \ \psi]^T$ reprezentálja az állapotvektort, mely a pályára merőleges sebességből és a legyezési szögsebességből áll, illetve $u = \delta$ a jármű kormányzójele.

A szabályzótervezés során az elsődleges követelmény a biztonságos mozgás, ebben az esetben ez a pályakövetési hiba minimalizálásaként van megfogalmazva. Szintén fontos követelmény a kormányzójele limitációja a beavatkozásszerv fizikai korlátait figyelembevéve, illetve esetleges oszcillációk kialakulását elkerülendő. Ezek miatt a szabályzótervezés a pályakövetési hiba és a kormányzójele normájának minimalizálásán alapult.

A \mathcal{H}_∞ szabályzótervezés során szükséges a megfelelő súlyozó függvények meghatározása a rendszert terhelő zajok, zavarások, bizonytalanságok és performanciák esetében, ezekre található részletes példa Sename et al. (2013) és Németh és Gáspár (2021) publikációikban.

A jármű oldalirányú irányítása mellett a robusztus hosszirányú irányítás is megtervezésre került. A jármű sebességét egy beépített PID kontrollert szabályozta, így a magas szintű szabályzás feladata az optimális sebességprofil meghatározása. Ahhoz, hogy a jármű biztosan követni tudja az adott pályát, a pályán való haladáshoz tartozó oldalirányú gyorsulásnak minden esetben a jármű által felvehető maximális alatt kell maradnia. Ez a határgyorsulás a tapadási körülményektől függ, az abroncs állapotától és az útviszonyoktól.

Amennyiben a pálya vonala előre ismert, úgy a biztonságos referencia sebesség a következőképpen határozható meg:

$$v_{ref}(s) \leq \eta \sqrt{\frac{a_{y,max}}{\kappa(s)}}, \quad (4)$$

ahol $a_{y,max}$ a maximálisan megengedhető oldalirányú gyorsulás a fellépő kerékerők alapján, κ a pályamenti görbületet leíró függvény, s a pálya paramétere, η pedig egy biztonsági tényező. A fenti (4)-es összefüggés a jármű tömegközéppontjának a tömegpont modellje alapján írható fel.

3.1 Megeősítéses tanulás alapú ágens tanítása

A megeősítéses tanulás alapú ágens egyszerre kezeli a jármű hosszirányú és oldalirányú dinamikáját egy magas minőségi irányítórendszerként. Nemlineáris volta képes kezelni a járműmodell nemlinearitásait, és a megadott jutalomfüggvény szerinti optimális irányítást hajtja végre.

Az ágens tanítása egy szimulált környezetbe történik, ami a jármű komplex nemlineáris modelljén alapul, figyelembe véve a nemlineáris gumikarakterisztikát, a kormánymű dinamikáját, a valós környezetben előforduló zajokat, zavarásokat, aktuátor szaturációt és időkésést. A ágens megvalósító neurális hálózat bemenetén a bejárandó pálya középvonalának pontjai találhatóak egyenletes osztásközzel. A feladat dimenziójának a csökkentése érdekében a pontok járműhöz viszonyított relatív értékei lettek figyelembevéve a jármű koordináta-rendszerében. Ezzel a megoldással a jármű mindig az origóban van és az x tengely pozitív irányába néz. A neurális háló bemeneti vektora tehát:

$$x_{NN} = [X_{r,1} Y_{r,1} X_{r,2} Y_{r,2} \dots X_{r,N} Y_{r,N}], \quad (5)$$

ahol N a számításnál felhasznált pályapontok számát jelzi. EZ egy tervezési paraméter, nagyobb értéke esetén több információ áll a háló rendelkezésére a jármű előtt lévő pályáról, azonban ezzel együtt a hálóparaméterek száma is növekszik, a tanítási folyamat hosszabb lesz. Túl kis értékek esetén pedig esetlegesen nem fog elég információ a hálózat rendelkezésére állni a manőverek végrehajtásához. Ezek miatt N értékének megválasztása főként tapasztalati úton történik a jellegzetes pályáívek geometriáját figyelembevéve.

A megeősítéses tanulás alapú ágens tanítása egy jutalomfüggvény alapján történik. Minden időlépésben a környezet állapotától függően ezen függvény alapján jutalmakat gyűjt az ágens, és a gyűjtött jutalmak összege jellemzi az ágens minőségét a tanítás során. A tanítás célja tehát a gyűjtött jutalom maximalizálása.

A feladat elvégzéséhez egy egyszerű paraméteres jutalomfüggvény került meghatározásra:

$$R(s_{env}, a) = -Ax_{Lat, err}^2 - B\delta_{ref}^2 + \Delta p, \quad (6)$$

ahol $x_{Lat, err}$ a pályakövetés oldalirányú hibája, Δp a jármű előrehaladása a középvonal mentén az adott időlépésben, δ_{ref} a szabályzó által választott referencia kormányaszög, s_{env} a környezet állapotát jelöli, a pedig az ágens által adott lépésben választott akciót. A tanítás során, abban az esetben, ha a jármű elhagyta a pálya széleit, akkor a (6)-on felül egy nagy negatív értékű jutalom kerül alkalmazásra, ezzel büntetve az ehhez hasonló viselkedését a rendszernek.

4. A BEMUTATOTT STRUKTÚRA ILLUSZTRÁCIÓJA

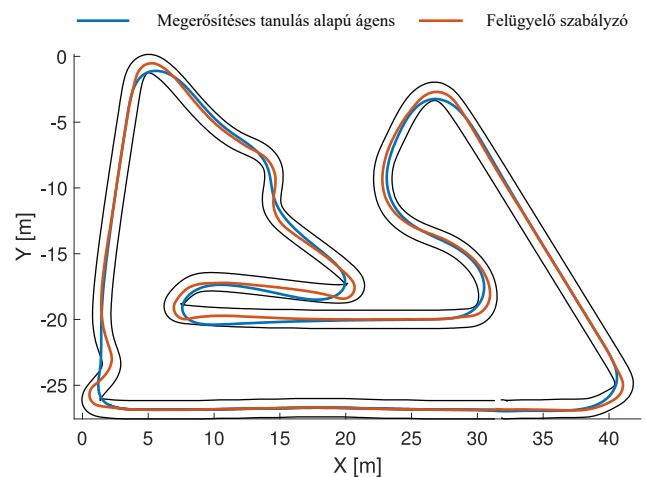
Ebben a fejezetben a bemutatott irányítási struktúra hatékonysága kerül illusztrálásra egy szimulációs példán keresztül. A szimuláció egy kisméretű tesztjármű modelljén alapszik, amely egy 1:10 arányban kicsinyített jármű (F'TENTH). A versenyautónak a Forma 1-es Bahreini Nemzetközi Versenypálya 1:30 arányban kicsinyített másán kell végig haladnia.

A felhasznált neurális hálózat 3 rejtett réteggel rendelkezett, melyek 16, 32 és 64 neuront tartalmaztak és ReLU típusú aktivációs réteget. A kimeneten hiperbolikus tangens aktivációs függvény került alkalmazásra, mely folytonosan deriválható módon képes figyelembe venni a kormány aktuátor korlátait és a jármű maximális sebességét. A háló bemenetén $N = 7$ db pályapont került figyelembevétele, melyek egyenletes 0,5 m-es távolságra helyezkednek el a pálya középvonalán, azaz az ágens 3,5 métert volt képes előre érzékelni a pályából. A maximális megengedett sebessége a járműnek $v_{max} = 3,5$ m/s. A jutalomfüggvény paramétereinek a numerikus értékei

$$R(s_{env}, a) = -0,4x_{Lat, err}^2 - 1,2\delta_{ref}^2 + \Delta p \quad (7)$$

összefüggés szerint adódtak, ami összességében a gyors előrehaladást, és a kormányaszög minimalizálását részesítette előnyben a pályakövetéssel szemben. Ezzel a cél a minél gyorsabb köridők teljesítése.

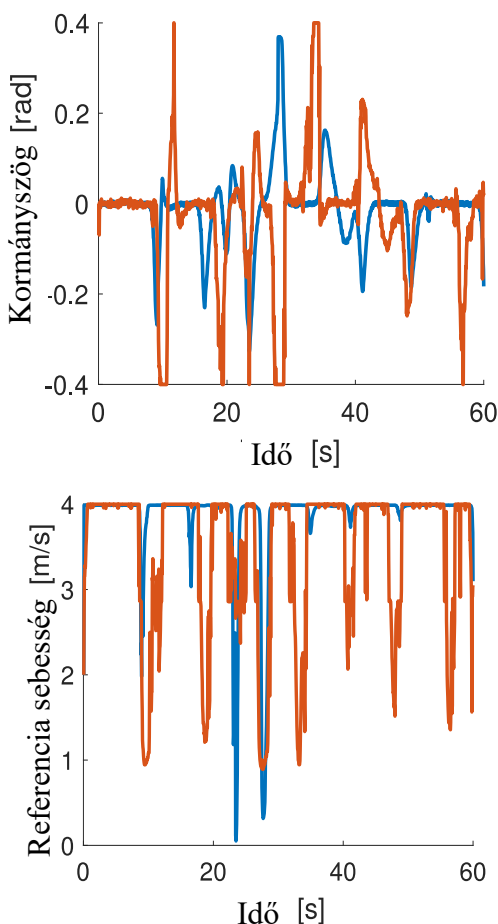
A felügyelő szabályzóban végrehajtott optimalizálási feladat során $\Delta L_\delta \in [-0,15; 0,15]$ és $\Delta L_v \in [-1; 0,1]$ intervallumok lettek meghatározva az a robusztus irányítójel környezetében. A felügyelőben található predikciós réteg egy kétlépéses becslést hajtott végre a jármű egyszerűsített modellje alapján a pályakövetési hiba korlátozása érdekében, a maximálisan megengedett eltérés $e_{max} = 0,4$ méter volt biztonsági okokból, hiszen ez jelentette a pálya két szélét.



2. ábra. A szimuláció során adódó trajektóriák a két szabályzó esetében.

A szimuláció során két szabályzó került összehasonlításra, a tisztán megerősítéses tanulást alkalmazó irányítási ágens, mely csupán a (7)-es egyenletben ismertetett jutalomfüggvény alapján adódó optimális irányítást hajtja végre, illetve a teljes felügyelő szabályzó alapú rendszer, mely az előbbi mellett tartalmazza a robusztus szabályzót is. A szimulációk eredményeként adódó trajektóriák láthatóak a 2. ábrán. Ezek alapján egyértelműen látszik, hogy a megerősítéses tanulás alapú ágens úgy volt képes gyorsabb köridőkre, hogy rendszeresen a pálya középvonalától való letérésre kényszerült, ezzel jobban közelítve a ideális ívet. Ezzel szemben a felügyelő szabályzót is tartalmazó rendszer lényegesen csökkentette pályakövetési hibát.

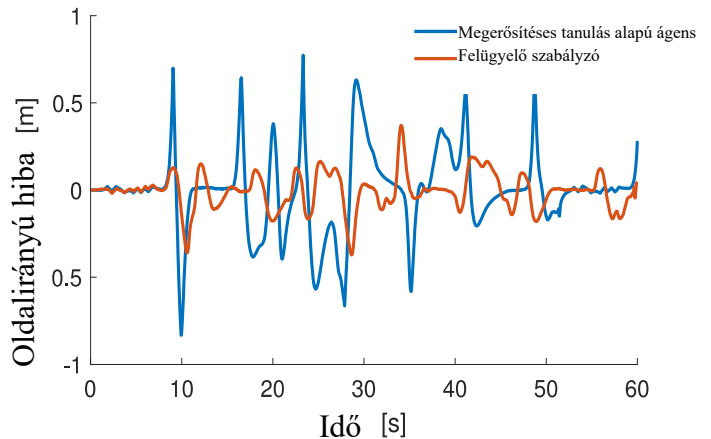
Hasonló következtetés vonható le a két rendszer által használt beavatkozó jelek tekintetében is. A 3. ábra mutatja a használt kormányzójelet és referenciasebességet. Látható, hogy a felügyelő szabályzó jóval konzervatívabban választja meg a referencia sebességet, azonban a pontosabb pályakövetések érdekében nagyobb kormánybeavatkozásokat is használ, ezzel csökkentve a felügyelt szabályzót alkalmazó struktúra (7) szerinti performanciáját. Mindennek az eredménye a köridő tekintetében 55 s a megerősítéses tanulás alapú ágens esetében és 60 s a felügyelő szabályzó alapú struktúra esetében. A



3. ábra. A szabályzók által felhasznált kormányzójelet és referencia sebesség.

biztonságos viselkedés ebben az esetben körülbelül 9%-os köridő növekedést eredményezett.

A fő eredménye a felügyelő alapú struktúrának ebben a példában a pályakövetési hiba sikeres korlátozása. Ennek illusztrálása látható a 4. ábrán, amely a szimuláció során adódó oldalirányú hibát szemlélteti. Egyértelműen látható, hogy a felügyelő szabályzót tartalmazó rendszer lényegesen csökkenti a pályakövetési hibát. Az előírt maximális érték $e_{max} = 0,4$ m volt, ami a megerősítéses tanulás alapú ágens esetében többször is megsértésre került, a legrosszabb esetben a 0,8 méter is elérte, míg a felügyelő szabályzó esetében



sikeresen teljesült ez a kritérium is.

4. ábra. A szimuláció során adódó pályakövetési hibák grafikonja.

Mindkét esetben elmondható, hogy a szabályzók képesek voltak a pályán történő navigálásra. A megerősítéses tanulás alapú ágens célja a köridő minimalizálása volt, ennek érdekében gyakran folyamodott a kanyarok levágásához a gyorsabb haladás, és a kisebb fékezések érdekében. A tanítás folyamat során nem volt szigorúan limitálva a pályakövetési hiba, csupán egy büntető tagként lett figyelembe véve a jutalom függvényben. Ezzel szemben a felügyelő szabályzóban történő kvadratikus optimalizálás és modell alapú predikció segítségével lehetőség adódott a pályakövetési hiba az előírt szint alatt tartására. Ezt a viselkedést a bemutatott szimuláció is alátámasztotta. A módszer hátránya, hogy korlátozások miatt nem képes olyan performanciára, mint a megerősítéses tanulás alapú ágens, ami a köridőkben tapasztalt különbségben nyilvánul meg, azonban a felügyelő szabályzó képest biztosítani a robusztus stabilitást a \mathcal{H}_∞ szabályzás révén, illetve a pályakövetési hiba limitációját a modell alapú optimalizáció révén.

5. KONKLÚZIÓ

Ez a tanulmány egy járműirányítási struktúra ismertetésével foglalkozott, amely pályakövető irányítások esetén használható nagy hatékonysággal. A struktúra három részből áll, egy klasszikus robusztus irányítási eljárásból, egy modern adat alapú nemlineáris szabályzásból és az ezek előnyeit egyésséítő felügyelő szabályzóból.

Az ismertetett struktúra hatékonysága bemutatásra került egy szimulációs példán keresztül, ahol a kapott eredmények egybeesnek az elvárta, a felügyelő szabályzó képes a biztonságos irányításra, a robusztus stabilitásra és a pályakövetési hiba limitációjára a performanica kismértékű csökkenése mellett.

A lehetséges jövőbeli kutatási irányok magába foglalják a módszer valós tesztjárművön történő implementációját és validációját. A tanítás során egy kisméretű tesztjármű modellje került beépítésre, mely tesztjármű rendelkezésre áll, és rendelkezik a szükséges érzékelőkkel és aktuátorokkal autonóm funkciók ellátásához. Az alkalmazhatóság szempontjából tehát rendkívül fontos a biztonságkritikus rendszerek valós környezetben történő validációja, melyet a felügyelő szabályzó alapú irányítórendszer esetében is szükséges kivitelezni.

KÖSZÖNETNYILVÁNÍTÁS

A publikációban szereplő kutatást a HUN-REN SZTAKI az Európai Unió támogatásával valósította meg, az Autonóm Rendszerek Nemzeti Laboratórium keretében. (RRF-2.3.1-21-2022-00002). A kutatást részben a Nemzeti Kutatási, Fejlesztési és Innovációs Hivatal támogatta az OTKA pályázat keretében (No. K 135512.).

Lelkó Attila munkáját az ÚNKP-23-3, a Nemzeti Kutatási, Fejlesztési és Innovációs Hivatal Új Nemzeti Kiválóság Programja támogatta.

REFERENCES

- Fliess, M. and Join, C. (2021). *Machine learning and control engineering: The model-free case*. In Arai, K., Kapoor, S., and Bhatia, R., editors, Proceedings of the Future Technologies Conference (FTC) 2020, Volume 1, pages 258–278, Cham. Springer International Publishing.
- Fényes, D., Hegedűs, T., Németh, B., Szabó, Z., and Gáspár, P. (2022). *Robust control design using ultra-local model-based approach for vehicle-oriented control problems*. In 2022 European Control Conference (ECC), pages 1746–1751.
- Haiyang, Z., Yu, S., Deyuan, L., and Hao, L. (2016). *Adaptive neural network pid controller design for temperature control in vacuum thermal tests*. In 2016 Chinese Control and Decision Conference (CCDC), pages 458–463.
- He, W., Yan, Z., Sun, Y., Ou, Y., and Sun, C. (2018). *Neural-learning-based control for a constrained robotic manipulator with flexible joints*. IEEE Transactions on Neural Networks and Learning Systems, 29(12):5993–6003.
- Huang, X., Kwiatkowska, M., Wang, S., and Wu, M. (2017). *Safety verification of deep neural networks*.
- Kabzan, J., Hewing, L., Liniger, A., and Zeilinger, M. N. (2019). *Learning- Based Model Predictive Control for Autonomous Racing*. IEEE Robotics and Automation Letters, 4(4):3363–3370.
- Khosravani, S., Khajepour, A., Fidan, B., Chen, S.-K., and Litkouhi, B. (2014). *Development of a robust vehicle control with driver in the loop*. In 2014 American Control Conference, pages 3482–3487.
- Kong, J., Pfeiffer, M., Schildbach, G., and Borrelli, F. (2015). *Kinematic and dynamic vehicle models for autonomous driving control design*. In 2015 IEEE Intelligent Vehicles Symposium (IV), pages 1094–1099.
- Lelkó, A., Németh, B., and Gáspár, P. (2021). *Stability and tracking performance analysis for control systems with feed-forward neural networks*. In 2021 European Control Conference (ECC), pages 1485–1490.
- McKinnon, C. D. and Schoellig, A. P. (2019). *Learn fast, forget slow: Safe predictive learning control for systems with unknown and changing dynamics performing repetitive tasks*. IEEE Robotics and Automation Letters, 4(2):2180–2187.
- Németh, B. and Gáspár, P. (2021). *Guaranteed Performances for Learning- Based Control Systems Using Robust Control Theory*, pages 109–142. Springer International Publishing, Cham.
- Rosolia, U. and Borrelli, F. (2020). *Learning How to Autonomously Race a Car: A Predictive Control Approach*. IEEE Transactions on Control Systems Technology, 28(6):2713–2719.
- Ruan, W., Huang, X., and Kwiatkowska, M. (2018). *Reachability analysis of deep neural networks with provable guarantees*.
- Senane, O. (2021). *Review on LPV approaches for suspension systems*. Electronics, 10(17):2120.
- Senane, O., Gáspár, P., and Bokor, J. (2013). *Robust Control and Linear Parameter Varying Approaches*. Springer Verlag, Berlin.
- Wu, M., Wicker, M., Ruan, W., Huang, X., and Kwiatkowska, M. (2020). *A game-based approximate verification of deep neural networks with provable guarantees*. Theoretical Computer Science, 807:298–329.
- Zhang, M., Wang, X., Yang, D., and Christensen, M. G. (2021). *Artificial neural network based identification of multi-operating-point impedance model*. IEEE Transactions on Power Electronics, 36(2):1231–1235.