



## FUZZY SZABÁLYBÁZIS OPTIMALIZÁLÁS A HFRIQ-LEARNING RENDSZERBEN

TOMPA TAMÁS

Miskolci Egyetem

Informatikai Intézet

Általános Informatikai Intézeti Tanszék

[tompa@iit.uni-miskolc.hu](mailto:tompa@iit.uni-miskolc.hu)

KOVÁCS SZILVESZTER

Miskolci Egyetem

Informatikai Intézet

Általános Informatikai Intézeti Tanszék

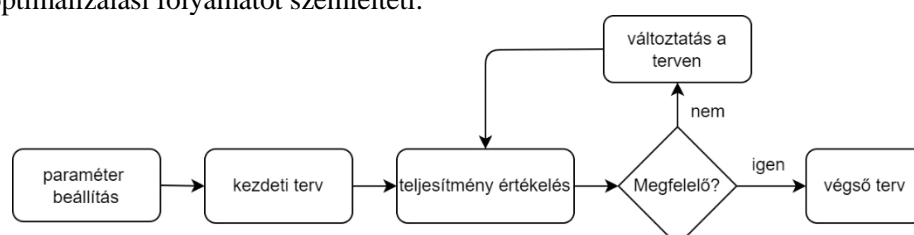
[szkovacs@iit.uni-miskolc.hu](mailto:szkovacs@iit.uni-miskolc.hu)

**Absztrakt.** A HFRIQ-learning (heurisztikusan gyorsított FRIQ-learning) a 'FIVE' fuzzy szabály-interpolációs módszeren alapuló Q-tanuló algoritmus, amelybe szakértő által megadott tudásbázis illeszthető fuzzy produkciós szabályok formájában. A tanulási folyamat során a kezdeti szakértői tudásbázist a rendszer optimalizálja olyan módon, hogy ha szükséges, akkor új fuzzy szabálypontot hoz létre, egyébként pedig a meglévő inkrementális szabályrendszert (Q-függvényt) hangolja. A hangolási folyamat során a szabálypontok pozíciója (antecedense és konzekvense) a gradiens módszer alkalmazása következtében a Q-függvény gradiense által módosul. A cikk célja annak bemutatása, hogy a tanulási folyamat során a szabálypontok optimalizálása hogyan valósul meg a HFRIQ-learning rendszerben.

*Kulcsszavak:* megerősítéses tanulás, heurisztikusan gyorsított megerősítéses tanulás, szakértői tudásbázis, Fuzzy szabály-interpoláció, Q-learning, FRIQ-learning

### 1. Bevezetés

Az optimalizálási módszerek célja változók olyan értékeinek megtalálása, amelyek minimalizálnak (vagy maximalizálnak) egy célfüggvényt. Ezen algoritmusok az adott változók értékeit addig módosítják (hangolják) amíg a valamilyen szempontból optimális megoldás elő nem áll. A következő ábra az általános optimalizálási folyamatot szemlélteti:

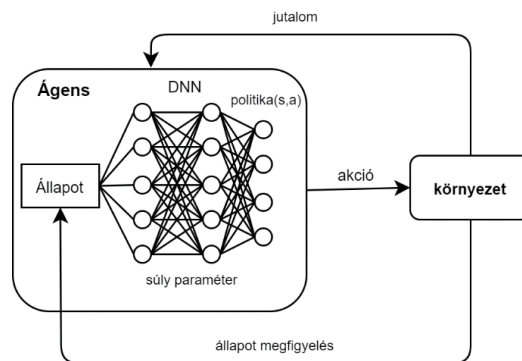


1. ábra: Általános optimalizálási folyamat [10]

Az optimalizálási folyamat általában egy függvény minimum vagy maximum pontjának (szélsőértékének) keresését jelenti. A legtöbb esetben egy hibafüggvény minimum pontjának keresése a cél, azaz a változók azon értékeinek a megtalálása, amely mellett a függvény értéke a legkisebb majd ez által úgy igyekszik elmozdítani (elhangolni) az adatpontokat, hogy a hiba értéke csökkenjen.

Az iteratív módszerek közül a gradiens módszer (gradient descent, GD – gradiens süllyedés) a legelterjedtebb, amely iterációkon keresztül jut el egy függvény minimumához, a függvény gradiense, azaz az iránymenti deriváltjai által úgy, hogy mindig a legmeredekebb lejtő irányba halad. Gradiens módszer alapú hangolási eljárást a mély tanulás [16] (Deep Learning) alapú „Deep Q-learning Network” (DQN) [26] is alkalmaz. A DQN egy olyan Q-learning algoritmus, amely a mély tanulás, azon belül is a mély megerősítéses tanulás (Deep Reinforcement Learning) [1][6][17] módszerek csoportjába sorolható. A mély megerősítéses tanulás a megerősítéses tanulás és a mély tanulás kombinációja. Ezen módszerek közös jellemzője, hogy általában nagyszámú bemeneti adattal dolgoznak és többretegű neurális hálózatokat (Deep Neural Network – DNN) alkalmaznak, melyek topológiájában több rejtett réteg is található (innen a „mély” elnevezés). Minden egyes réteg egy másfajta reprezentációját adja a bemeneti adatoknak, amely egyfajta jellemzőkinyerésnek tekinthető, így rétegről-rétegre előre haladva egyre komplexebb összefüggések kerülnek beazonosításra (például: pixelek → élek → alakzatok → tárgyak). Az elterjedtebb mély megerősítéses tanulási módszerekről bővebb áttekintést a [1][6] és [17] hivatkozások adnak. Az optimalizálási módszereket a Deep Learning algoritmusok esetében általában a neurális hálózat súlyainak optimalizálására alkalmazzák, amely által az úgy igyekszik módosítani a hálózat súlyainak értékét, hogy a hiba értéke csökkenjen.

A mély megerősítéses tanulás modelljét a következő ábra szemlélteti [16]:



2. ábra: A mély megerősítéses tanulás modellje [16]

A DQN módszer neurális hálózatot alkalmaz a Q-függvény leírására, amely által a Q-függvény a következő módon írható fel [4]:

$$Q(s, a) \approx Q(s, a, \theta) \quad (1)$$

ahol  $Q(s, a)$  a Q-függvény,  $s$  az állapot,  $a$  az akció, a  $\theta$  pedig a neurális hálózat súlyait reprezentálja. A hiba ( $L$ ) ebben az esetben a TD-hiba alapján határozható meg, amelyet az adott optimalizálási módszer minimalizálni igyekszik [4]:

$$L = E \left[ \left( r_k + \gamma \max_a Q(s_{k+1}, a_{k+1}, \theta) - Q(s_k, a_k, \theta) \right)^2 \right] \quad (2)$$

Ahol  $r_k + \gamma \max_a Q(s_{k+1}, a_{k+1}, \theta) - Q(s_k, a_k, \theta)$  összefüggés a TD-hiba,  $k$  az iteráció (időpillanat),  $r_k$  a jutalom,  $\gamma$  diszkontálási tényező,  $s_k$  az állapot,  $a_k$  az akció,  $\theta$  pedig a neurális hálózat súlyai.

A klasszikus gradiens módszernek több változata is található a szakirodalomban, ilyen például a sztochasztikus gradiens módszer (Stochastic Gradient Descent - SGD), amely a klasszikus gradiens módszerrel ellentétben nem az összes lehetséges hibafüggvény pontban határozza meg a gradienst és ez által az összes rendelkezésre álló minta alapján hangol minden egyes iterációban, hanem egyesével, véletlenszerűen veszi a mintapontokat iterációnként. Ennek előnye, hogy nagy dimenziószámmal rendelkező problémák esetében csökken a számítási- és memória igény. A „mini-batch” gradiens módszer a mintaadatok halmazát kisebb részhalmazokra bontja, majd ezen részhalmazok alapján határozza meg a hibát és frissíti (azaz hangolja) a modell paramétereit. A gradiens módszer alapú optimalizálási eljárások hátránya, hogy nem garantálják a globális optimum megtalálását. A gradiens módszerek elterjedtebb változatairól és azok hatékonyságának összehasonlításáról a [7] és [19] irodalmak adnak bővebb áttekintést.

A részecske-raj alapú optimalizálás (Particle Swarm Optimization - PSO) [8] szintén egy iteratív algoritmus, amely működése a keresési térben, általában egyenletes eloszlás szerint elhelyezett részecskéken (raj) alapszik, amely részecskék matematikai összefüggések alapján mozognak. A részecskék a legjobb pontot keresik a térben majd a saját legjobb pozíciójuk és a raj legjobb pozíciója alapján mozdulnak el. A legjobb ismert pozíció frissül iterációnként attól függően, hogy a raj talált-e a legjobb ismert pozíciónál még jobbat vagy sem, tehát a részecskék mozgását (és a mozgás irányát) az adott iterációban megtalált legjobb pozíció befolyásolja. Ha a raj már nem talál jobb pozíciót, akkor a leállási feltétel bekövetkezése után a raj legjobbjá (részecskéje) adja meg a függvény optimum, azaz minimum pontját. A PSO algoritmus esetében nincs szükség a függvény gradiensének és így annak parciális deriváltjainak meghatározására (a gradiens módszereknél ez alapkövetelmény) így ez a módszer jól alkalmazható olyan függvények esetében melyek gradiense nem ismert illetve zajjal terhelt, valamint időben változó problémák esetében is.

További optimalizálási módszerekről a [10] sorszámú szakirodalom, a megerősítéses tanulásban alkalmazott elterjedtebb optimalizálási módszerekről pedig a [18] hivatkozás ad bővebb áttekintést.

## 2. Heurisztikusan gyorsított FRIQ-learning (HFRIQ-learning)

A heurisztikusan gyorsított FRIQ-learning (Heuristically Accelerated Fuzzy Rule-Interpolation based Q-learning - HFRIQ-learning) [24] a FRIQ-learning (Fuzzy Rule-Interpolation based Q-learning) [28] módszer kiterjesztése, amely által az alkalmas külső szakértői tudásbázis beágyazására [20] illetve hangolására [21][24]. A módszer a „FIVE” (Fuzzy Rule Interpolation based on Vague Environment) [109] fuzzy szabály-interpolációs eljárást alkalmazza a Q-függvény leírására, amely által az állapot-akció tér folytonos.

### 2.1. A „FIVE” Fuzzy szabály-interpolációs módszer

A „FIVE” Kovács-Kóczy [11][12][14] által kifejlesztett egylépéses szabály-interpolációs módszer, amely az interpolációs feladatot egy úgynevezett bizonytalan környezetbe (Vague Environment, VE) [9] helyezi át. A Klawonn-féle bizonytalan környezet alapgondolata az univerzum elemei közötti hasonlóságon illetve megkülönböztetlenségen alapszik [9]. Az elemek hasonlóságának mértéke súlyozott távolság által definiálható, ahol a súlytényező az  $s(x)$  skálafüggvény [9][12][13]:

$$s(x) = |\mu'(x)| = \left| \frac{d\mu}{dx} \right| \quad (3)$$

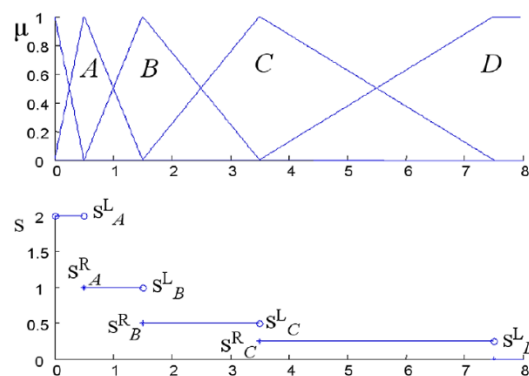
$$\text{létezik ha: } \min\{\mu_i(x), \mu_j(x)\} > 0 \rightarrow |\mu'_i(x)| = |\mu'_j(x)| \quad (4)$$

Ahol  $\forall i, j \in I$  fuzzy halmazok,  $\mu'(x)$  a fuzzy halmaz tagsági függvényének deriváltja, ebben az esetben tehát az  $s(x)$  skálafüggvény a  $\mu(x)$  tagsági függvény deriváltja, [12]. Az  $(x_1, x_2)$  elemek közötti, bizonytalan környezetbeli  $\delta_s$  távolságuk meghatározásának módja az  $s(x)$  skálafüggvényen alapszik [9]:

$$\delta_s(x_1, x_2) = \left| \int_{x_2}^{x_1} s(x) dx \right| \quad (5)$$

Az  $X$  bizonytalan környezetben az  $(x_1, x_2)$  elemek  $\varepsilon$  mértékben hasonlóak (megkülönböztetetlenek) tekinthetők, ha a közöttük lévő  $\delta_s$  távolság legfeljebb  $\varepsilon$  mértékű, azaz  $\delta_s(x_1, x_2) \leq \varepsilon$  [12][13]. A bizonytalan környezetek (antecedens, konzekvens, szabálybázis) előre számolhatók (így biztosítva a módszer gyorsaságát) amelyben minden szabály egy-egy szabálypontként ábrázolható.

Az alábbi 3. ábra fuzzy halmazokat (az ábra felső részében lévő grafikon) és az őket jellemező skálafüggvényeket (az ábra alsó részében lévő grafikon) szemlélteti, amelyek ezáltal alkalmasak az adott fuzzy partíció alakjának leírására [13]:



3. ábra: Fuzzy halmazok (felső grafikon) és az őket jellemező skálafüggvények (alsó grafikon) [13]

Minél kisebb a skálafüggvény értéke az elemek annál kevésbé megkülönböztethetők egymástól, azaz ugyanolyan alaphalmazbeli távolság esetében egyre közelebb vannak egymáshoz. Ha a skálafüggvény értéke nulla az azt eredményezi, hogy az elemek nem megkülönböztethetők, mert egyformán közel helyezkednek el. Egy valós gyakorlati alkalmazás esetében ez azt eredményezi, hogy például minden olyan esetben amikor 2 méternél messzebb találhatók objektumok akkor nem kell

fékezni, amikor pedig ennél közelebb vannak akkor egyre jobban kell fékezni.

A „FIVE” módszer a multidimenziós mivolta következtében a Shepard interpolációs operátort [5] alkalmazza, amely által a singleton (egyértékű)  $c_k$  konzekvens a következő összefüggés alapján, további defuzzifikációs lépések nélkül határozható meg [13]:

$$y(\mathbf{x}) = \begin{cases} c_k & \text{ha } \mathbf{x} = \mathbf{a}_k \text{ minden } k\text{-ra,} \\ \sum_{k=1}^r \left( \left( \frac{c_k}{\delta_{s,k}^\lambda} \right) / \left( \sum_{k=1}^r \frac{1}{\delta_{s,k}^\lambda} \right) \right) & \text{egyébként} \end{cases} \quad (6)$$

Ahol  $\mathbf{x}$  a többdimenziós megfigyelés,  $c_k$  a  $k$ -adik létező szabály konzekvens,  $r$  a szabályok száma az  $R$  szabálybázisban,  $\lambda$  a Shepard-kivető,  $\delta_{s,k}$  pedig a súlyozott (Euklideszi) távolság, amely a következő formával által definiálható [12][13]:

$$\delta_{s,k} = \delta_s(\mathbf{a}_k, \mathbf{x}) = \left[ \sum_{i=1}^m \left( \int_{a_{k,i}}^{x_i} s_{x_i}(x_i) dx_i \right)^2 \right]^{1/2} \quad (7)$$

Ahol  $\mathbf{x}$  az  $m$ -dimenziós crisp megfigyelés,  $\mathbf{a}_k$  a magja az  $m$ -dimenziós szabály antecedens  $\mathbf{A}_k$ -nak,  $s_{x_i}$  pedig az  $i$ -edik skálafüggvény az  $m$ -dimenziós antecedens univerzumban.

A „FIVE” [13] tehát egy alkalmazás-orientált FRI (Fuzzy Rule Interpolation) módszer, amely alacsony számítási igénye miatt [2][3] jól használható valós idejű alkalmazásokban illetve robotikai irányításokban. Továbbá, ezen alkalmazott FRI módszer által a rendszer komplexitása csökkenthető a ritka szabálybázis következtében illetve a rendszer abban az esetben is szolgáltat kimentet mikor a klasszikus fuzzy következtetési eljárások nem [15].

## 2.2. A HFRIQ-learning

A HFRIQ-learning rendszer tudásbázisa egy ritka fuzzy szabálybázis által leírt, egy  $r_i$  ( $i \in [1, m]$ ) szabály formája az  $m$  méretű  $R$  szabálybázisban a következő [28]:

$$r_i: \text{If } s_1 \text{ is } S_1^i \text{ And } s_2 \text{ is } S_2^i \text{ And ... And } s_n \text{ is } S_n^i \text{ And } a \text{ is } A^i \text{ Then } \tilde{Q}(s, a) = q^i \quad (8)$$

ahol  $S_j^i$  az  $i$ -edik ( $i \in [1, m]$ ) szabály  $j$ -edik ( $j \in [1, n]$ ) állapot dimenziójának fuzzy halmaza az  $n$ -dimenziós  $\mathcal{S}$  állapottérben,  $\mathbf{s} \in \mathcal{S}$  az  $n$ -dimenziós állapot megfigyelés,  $s_j$  a  $j$ -edik dimenziója az  $\mathbf{s}$  állapot megfigyelésnek,  $A^i$  az  $i$ -edik szabály egydimenziós akció univerzumának ( $U$ ) fuzzy halmaza,  $a \in U$  az akció,  $\tilde{Q}(s, a)$  a FIVE FRI [13] által becsült Q-függvény,  $q^i$  pedig az  $i$ -edik szabály konzekvens (Q-értéke).

Az  $R_{expert}$  szakértői tudásbázis formátuma hasonló az (8) formula által definiált fuzzy szabályokhoz, azzal az eltéréssel, hogy az  $\hat{r}$  szakértői szabályok antecedense az állapot, konzekvens pedig az ebben az állapotban preferált akció [20]:

$$\hat{r}_i: \text{If } s_1 \text{ is } \hat{S}_1^i \text{ And } s_2 \text{ is } \hat{S}_2^i \text{ And ... And } s_n \text{ is } \hat{S}_n^i \text{ Then } a = \hat{A}^i \quad (9)$$

ahol  $\hat{r}_i$  az  $i$ -edik ( $i \in [1, \hat{m}]$ ) szakértői szabály az  $R_{expert}$  szabálybázisban,

$\hat{S}_n^i = [\hat{S}_1^i, \hat{S}_2^i, \dots, \hat{S}_n^i]$  az  $i$ -edik szakértői szabály  $n$ -dimenziós állapot megfigyelése,  $\hat{A}^i$  az ehhez az  $\hat{S}_n^i$  állapot megfigyeléshez tartozó akció,  $i$  ( $i \in [1, \hat{m}]$ ) pedig a szabály indexe az  $\hat{m}$  méretű szakértői szabályrendszerben.

Annak következtében, hogy a szakértői szabályrendszer injektálható legyen a rendszerbe szükséges a formátumának átalakítása állapot-akció-Q-érték formátumra. Ekkor az átalakított szakértői szabályok antecedense az állapot-akció, konzekvensze pedig egy becsült  $\tilde{Q}_{init}$  érték lesz. A becsült kezdeti  $\tilde{Q}_{init}$  érték a környezet által maximálisan adható megerősítés ( $g_{max}$ ) ismeretében határozható meg [20]. A HFRIQ-learning rendszer tanulási folyamata ezen  $R_{expert}$  szakértői szabályrendszer és a  $2^{n+1}$  darabszámú ( $n$ : állapotdimenziók száma), 0 konzekvens értékkel rendelkező  $r_i^{\square}$  sarokponti szabályok összefésülésével létrejött fuzzy szabálybázissal kezdődik [20][28]. Abban az esetben ha ellentmondás alakul ki, azaz szakértői szabály sarokponti szabályra illeszkedik (de eltérő a konzekvensük), akkor az ellentmondás feloldása következtében ezen két szabály összevonásra kerül egyetlen szabállyá. A létrejött kezdeti szakértői szabályrendszer a tanulási folyamat során inkrementálisan növekszik a rendszer által létrehozott szabályokkal [29][28]. Új szabály akkor kerül beillesztésre a szabálybázisba, ha a Q-frissítés ( $\Delta\tilde{Q}$ ) értéke nagyobb, mint egy  $\varepsilon_Q$  érték ( $\Delta\tilde{Q} > \varepsilon_Q$ ) [29] és a legközelebbi szabálypont is távolinak tekinthető [23][24]. A szabályközelség meghatározásának az alapja a szabályok között definiált, dimenzióként számított távolságok [23][24]. Abban az esetben, ha a  $\Delta\tilde{Q}$  érték kicsi ( $\Delta\tilde{Q} < \varepsilon_Q$ ), akkor a teljes szabálybázis konzekvensze kerül frissítésre a következő módon [29]:

$$\tilde{Q}^{k+1}(s, a) = \tilde{Q}^k(s, a) + \Delta\tilde{Q}^{k+1}(s, a) \quad (10)$$

$$\Delta\tilde{Q}^{k+1}(s, a) = \alpha * \left( g(s, a, s') + \gamma * \max_{a' \in U} \tilde{Q}^k(s', a') - \tilde{Q}^k(s, a) \right) \quad (11)$$

ahol  $\gamma \in [0,1]$  a leszámítolási tényező,  $\alpha \in [0,1]$  a tanulási ráta,  $a$  pedig az  $s$ -ben végrehajtott akció. Az új megfigyelt állapot  $s'$ ,  $g(s, a, s')$  a megfigyelt jutalom az  $s \rightarrow s'$  állapot átmenetre,  $\tilde{Q}^k$  és  $\tilde{Q}^{k+1}$  pedig a  $k$ -edik és a  $(k+1)$ -edik iteráció FIVE FRI módszer által becsült Q-értéke [29]:

$$\tilde{Q}(s, a) = \begin{cases} \sum_{i=1}^m \left( \left( \frac{q^i}{(\delta_v^i)^\lambda} \right) / \left( \sum_{j=1}^m \frac{1}{(\delta_v^j)^\lambda} \right) \right) & \text{ha } (s, a) = (s^i, a^i) \\ & \text{valamennyi } i - re, \\ & \text{egyébként} \end{cases} \quad (12)$$

ahol  $q^i$  az  $i$ -edik ( $i \in [1, m]$ ) szabály konklúziója,  $(s, a)$  a megfigyelés,  $\delta_v^i$  a skálázott távolság [13] az  $(s, a)$  megfigyelés és az  $i$ -edik szabály  $(s^i, a^i)$  antecedense között,  $\lambda$  a Shepard paraméter,  $m$  pedig a szabályok száma.

Ha a  $\Delta\tilde{Q}$  értéke kicsi és van már létező szabály a megfigyelés közelében, akkor a gradiens módszer alapú hangolási eljárás a megfigyeléshez legközelebb elhelyezkedő szabálypont antecedensét és konzekvensét fogja hangolni [24]. A szabálypont új pozíciója az alábbi módon kerül meghatározásra [24]:

$$s_{k+1} = s_k - \left( 2 * TDerror * \frac{\partial \tilde{Q}(s, a)}{\partial s} \right) * \alpha \quad (13)$$

$$a_{k+1} = a_k - \left( 2 * TDerror * \frac{\partial \tilde{Q}(s, a)}{\partial a} \right) * \alpha \quad (14)$$

$$q_{k+1} = q_k - \left( 2 * TDerror * \frac{\partial \tilde{Q}(s, a)}{\partial q} \right) * \alpha \quad (15)$$

ahol a  $s_{k+1}, a_{k+1}, q_{k+1}$  a gradiens-módszer által meghatározott új állapot, akció és Q-értékek,  $s_k, a_k, q_k$  a régi állapot, akció és Q-értékek,  $\alpha$  a gradiens-módszer tanulási rátája,  $\frac{\partial \tilde{Q}(s, a)}{\partial s}, \frac{\partial \tilde{Q}(s, a)}{\partial a}, \frac{\partial \tilde{Q}(s, a)}{\partial q}$  a Q-függvény állapot, akció és Q-érték szerinti parciális deriváltjai, a  $TDerror$  értéke pedig a következő [24]:

$$TDerror = g(s, a, s') + \gamma * \max_{a' \in \tilde{U}} \tilde{Q}^k(s', a') - \tilde{Q}^k(s, a) \quad (16)$$

Az alkalmazott hangolási módszer következtében, a szabálypontok vándorlása miatt előfordulhat olyan eset, hogy több szabálypont is közel kerül egymáshoz. Ebben az esetben az egymáshoz közel kerülő és ez által hasonló információt leíró szabálypontok egyesítésre kerülnek egyetlen szabállyá [23][24], ami által a szabálybázis mérete a tanulási folyamat csökkenthető. A szabálytávolság alapú szabálybázis redukálási módszerről bővebb információ a [23][24][25] hivatkozásokban található. További, a tanulási folyamat után opcionálisan alkalmazható szabálybázis csökkentési módszereket a [22][27][30] hivatkozások mutatnak be.

A HFRIQ-learning tanulási folyamata akkor ér véget, ha nem kerül új szabály hozzáadásra az inkrementális szabálybázisba, a Q-frissítés értéke elenyészően kicsi, a létező szabálypontok pozíciója nem változik, és nem kerülnek szabályok összevonásra.

### 3. Fuzzy szabálypontok hangolása a HFRIQ-learning rendszerben

A szakértő által definiált a priori szabályrendszerben amennyiben a megadott szakértői produkciós szabály valamelyik állapothoz nem megfelelő akció értéket rendel (azaz a szakértő szabályrendszer csak részben tekinthető helyesnek), úgy a szabálybázisba felvett szakértői szabály állapot-akció antecedense is rossz helyre kerül. Ebben az esetben, a csak részben helyes szabályrendszer negatív hatással lehet a tanulási folyamat hatékonyságára [20][21], így szükség lehet egy hangolási eljárásra, amely képes a szabályok állapot-akció pontját, azaz antecedensét elhangolni (optimalizálni) a megfelelő irányba, szabálypontba.

#### 3.1. Hangolandó szabálypontok meghatározása

Az FRIQ-learning Q-függvény tartópontjainak hangolása (optimalizálása) során meghatározandó, hogy mely szabálypontok kerüljenek optimalizálásra a tanulási fázis során. A módszerben [20][24] a Q-függvény a fuzzy szabálypontok „FIVE” FRI módszer fuzzy interpolációjával állnak elő.

A gradiens módszernek több verziója ismert a bemutatottak alapján. A GD módszer minden egyes iterációban minden egyes mintapontot figyelembe vesz a hibafüggvény illetve a deriváltak számításakor, míg a SGD módszer egyesével (véletlenszerűen) veszi a mintapontokat majd ezek alapján határozza meg a gradienst és hibafüggvényt iterációként. Jelen esetben mintapontoknak a

szabálypontok tekinthetők, melyek hangolása a TD-hiba (16) értékének függvényében történik.

A fuzzy szabályok antecedens és konzekvens értékeinek hangolása a  $x_{k+1} = x_k - \nabla F(x_k) * \alpha$  összefüggés alapján történik, ahol az  $F(x_k)$  függvény parciális deriváltja ( $\nabla F(x_k)$  gradiens) a láncszabály alkalmazásával a következőképpen határozható meg:

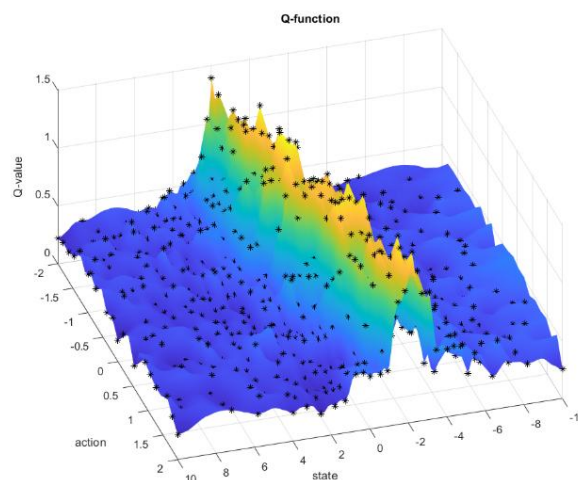
$$\nabla F(x_k) = \frac{\partial MSE(x_k)}{\partial x_k} = \frac{\partial (TDerror)^2}{\partial x_k} = 2 * TDerror * \frac{\partial \bar{Q}^k(s, a)}{\partial x_k} \quad (17)$$

A tanulási folyamat során a szabálypontok abban az esetben hangolódnak, ha a Q-frissítés értéke nagynak tekinthető ( $\Delta \bar{Q} > \varepsilon_Q$ ) és az éppen aktuális állapot-akció megfigyelés közelében (meghatározott közelségmérték alapján [23][24][25]) található már létező szabály.

Annak vizsgálata, hogy az összes szabálypont antecedens és konzekvens értékeinek hangolása a tanulási folyamat során milyen hatással van a Q-függvényre egy egyszerű, egy állapot- és egy akciódimenzióval rendelkező mintapéldán történik. A mintapélda paraméterei a következők:

- állapotváltozó:  $s_1 \in [-10, 10]$
- akcióváltozó:  $a \in [-2, 2]$
- $\alpha = 0.5$
- $\gamma = 0.4$
- $\varepsilon = 0.5$
- $g_{max} = 1$
- szakértői szabály: 0 állapotpontban 0 értékű akció
- jutalomfüggvény: a környezet +1 jutalmat ad, ha az ágens állapotváltozójának értéke -1 és +1 közötti, ellenkező esetben a jutalom értéke 0

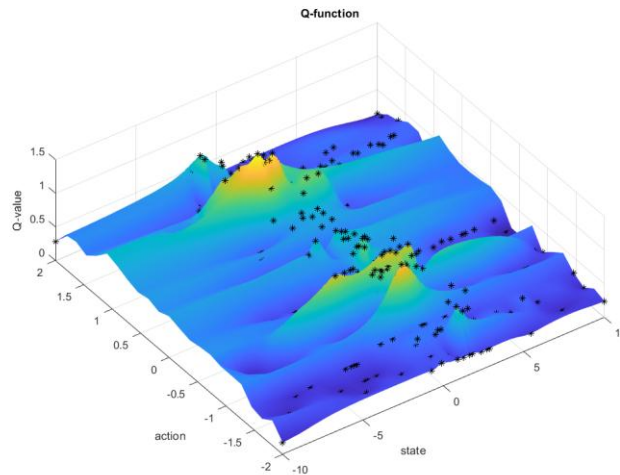
Az 4. ábrán az összehasonlítás alapjául szolgáló, a gradiens módszeren alapuló hangolási eljárás [24] alkalmazása nélkül kapott „referencia” Q-függvény felülete látható, ahol a „\*” (csillag) karakterek a szabálypontokat (530 darab) jelölik, a szabályok között megengedett minimális távolság pedig az univerzumok hosszának 100-ad része:



4. ábra: Az összehasonlítás alapjául szolgáló "referencia" Q-függvény felülete



Az első vizsgált esetben a javasolt hangolási módszerrel [24] az összes szabálypont hangolásra kerül, beleértve a szakértői szabályokat is. Az ebben a futási esetben kapott Q-függvény felületét a 5. ábra szemlélteti:



5. ábra: Összes létező szabálypont hangolása esetében kapott Q-függvény felülete

Megállapítható, hogy az összes szabálypont hangolásával a 4. ábrán látható referencia függvény felületéhez képest a függvény felülete romlott, nem alakult ki a felület „gerince”, több kisebb csúcs keletkezett csupán. Ennek oka az lehet, hogy az állapot-akció tér azon pontjainak antecedense és Q-értéke, melyek ritkán kerülnek bejárásra és így frissítésre (a felderítés-kiaknázás következtében), az összes szabálypont hangolása miatt elromlik. A kiaknázás során bejárt út frissítései hatással vannak a ritkán, csak a felderítés során tesztelt területekre. A visszajelzés nélküli (kevésbé bejárt) területek Q-értéke a gyakorta bejárt területek átlag Q-értéke felé hangolódik. Emiatt egyes szabályok Q következmény értéke elromlik, azaz rossz irányban módosul.

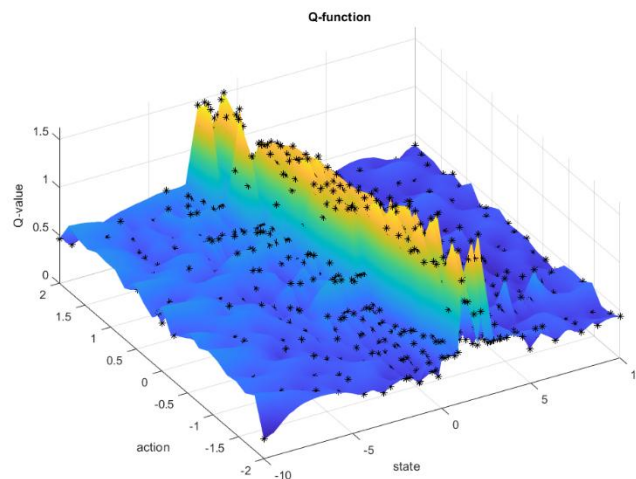
Az összes szabálypont egyidejű hangolása ezért nem járható út. Megoldás lehet az, hogy csak azon szabályok kerüljenek hangolásra, amely az éppen aktuális állapot-akció megfigyelési pont közelében található, azaz a legközelebb helyezkednek el ahhoz. A javasolt módszer megoldja a fentebb említett problémát és helyes irányban hangolja a Q-függvény szabálypontjait, nem változtatva azon szabályokat melyek a aktuális állapot-akció megfigyelési ponttól távol, a kevésbé látogatott területekre esnek. A javasolt módszerben a közelség mértékének meghatározása a [23][24][25] publikációkban bemutatottak alapján történik.

A hangolási folyamat során tekintettel kell lenni a hangolandó (az éppen közeli) szabálypont típusára. A sarokponti típusú szabályok antecedensei kötöttek (az állapot-akció tér univerzum hiperkocka sarokpontjai), a hangolás során nem változhatnak. A hangolás ezért a sarokponti típusú szabályok esetében csak a szabályok konzekvensére vonatkozik. A szakértői és a rendszer által felvett típusú szabályok esetén a szabályok antecedense (állapot-akció pontja) és konzekvensé (Q-értéke) is hangolásra kerül.

Abban az esetben, ha a tanulási folyamat során valamely szakértői szabály antecedense kerül hangolására, akkor ezen produkciós (állapot-akció) szakértői szabály módosul. Abban az esetben, ha a hangolási folyamat során az antecedens nem, vagy csak kismértékben változik, akkor a szakértői szabály a környezeti tesztek által igazoltan, helyesen definiált szakértői szabálynak tekinthető.

A 6. ábrán az az eset látható mikor nem az összes szabálypont, hanem csak a

megfigyeléshez közeli szabályok kerültek hangolásra. Ebben az esetben a függvény felülete nem romlott el, a függvény gerince viszonylag összefüggően kialakult, a függvény formája a referencia Q-függvényhez hasonló lett, azzal a különbséggel, hogy az kisimult, a felülete simább lett (a gradiens alapú hangolásnak köszönhetően).



6. ábra: Csak a megfigyeléshez legközelebbi szabálypontok hangolása esetében kapott Q-függvény felülete

A kapott Q-függvény felületeken sok olyan szabálypont található (a „\*” karakter által jelölt pozíciókban), melyek viszonylag egymáshoz közel (és sűrűn) helyezkednek el. Ezeken az állapot-akció területeken a gradiens módszer alapú hangolási eljárásnak köszönhetően a közeli szabályok egyre közelebb kerülnek egymáshoz. A hangolási folyamat következő lépése az, hogy az egymáshoz viszonylag közel kerülő szabályok, a szabályok közötti távolság, illetve távolságkülbszöbök alapján kerülnek összevonásra (egyesítésre), a szabályok számának csökkentése érdekében [23][24][25].

#### 4. Köszönetnyilvánítás

„A Kulturális és Innovációs Minisztérium ÚNKP-23-4-I-ME/5 kódszámú Új Nemzeti Kiválóság Programjának a Nemzeti Kutatási, Fejlesztési és Innovációs Alapból finanszírozott szakmai támogatásával készült.”



#### 5. Összefoglalás

Bemutatásra került a HFRIQ-learning rendszer gradiens módszer alapú szabálybázis optimalizálási módszere. Mintapélá segítségével igalásra került, hogy az összes szabálypont egyidejű hangolása negatív hatással van a Q-függvényre annak következtében, hogy a kiaknázás során bejárt út frissítései hatással vannak a ritkán,

csak a felderítés során tesztelt területekre. Így a kevésbé bejárt területek Q-értéke a gyakorta bejárt területek átlag Q-értéke felé hangolódik, ami miatt az egyes szabályok Q következmény értéke elromlik, azaz rossz irányban módosul. Bemutatásra került egy javasolt módszer, amely ezt a problémát kiküszöböli. A bemutatott módszer következtében a HFRIQ-learning megerősítéses tanulási rendszer tudásbázisának hangolása során az állapot-akció tér ritkán bejárt területein lévő fuzzy Q-szabályok elhangolódása csökkenthető, ha az összes fuzzy szabálypont egyidejű hangolása helyett, csak azon szabályok kerülnek hangolásra, amely az éppen aktuális állapot-akció megfigyelési pont közelében található.

### Irodalomjegyzék

- [1] Arulkumaran, Kai, et al. "Deep reinforcement learning: A brief survey." *IEEE Signal Processing Magazine* 34.6 (2017): 26-38. <https://doi.org/10.1109/msp.2017.2743240>
- [2] Bartók, Roland, and József Vásárhelyi. "Design of a FPGA accelerator for the FIVE fuzzy interpolation method." *International Journal of Computer Applications in Technology* 68.4 (2022): 321-331. <https://doi.org/10.1504/ijcat.2022.125185>
- [3] Bartók, Roland, and József Vásárhelyi. "Examining Cache Handling of the FIVE Method on Multicore Systems." 2019 IEEE 17th World Symposium on Applied Machine Intelligence and Informatics (SAMI). IEEE, 2019. <https://doi.org/10.1109/sami.2019.8782721>
- [4] Brunton, Steven L., and J. Nathan Kutz. *Data-driven science and engineering: Machine learning, dynamical systems, and control*. Cambridge University Press, 2019. <https://doi.org/10.1017/9781009089517>
- [5] D. Shepard, "A two dimensional interpolation function for irregularly spaced data", *Proc. 23rd ACM Internat. Conf.*, 1968, pp. 517-524. <https://doi.org/10.1145/800186.810616>
- [6] François-Lavet, Vincent, et al. "An introduction to deep reinforcement learning." *arXiv preprint arXiv:1811.12560* (2018). <https://doi.org/10.1561/9781680835397>
- [7] Haji, Saad Hikmat, and Adnan Mohsin Abdulazeez. "Comparison of optimization techniques based on gradient descent algorithm: A review." *PalArch's Journal of Archaeology of Egypt/Egyptology* 18.4 (2021): 2715-2743.
- [8] Kennedy, James, and Russell Eberhart. "Particle swarm optimization." *Proceedings of ICNN'95-international conference on neural networks*. Vol. 4. IEEE, 1995. <https://doi.org/10.1109/icnn.1995.488968>
- [9] Klawonn, F.: *Fuzzy Sets and Vague Environments*, in *Fuzzy Sets and Systems*, Vol. 66, 1994, pp. 207-221. [https://doi.org/10.1016/0165-0114\(94\)90311-5](https://doi.org/10.1016/0165-0114(94)90311-5)
- [10] Kochenderfer, Mykel J., and Tim A. Wheeler. *Algorithms for optimization*. MIT Press, 2019.
- [11] Kovács, Sz., Kóczy, L. T.: *Approximate Fuzzy Reasoning Based on Interpolation in the Vague Environment of the Fuzzy Rule base as a Practical Alternative of the Classical CRI*. *Proceedings of the 7th International Fuzzy Systems Association World Congress, Prague, Czech Republic, 1997*, pp. 144-149.
- [12] Kovács, Sz., Kóczy, L. T.: *The use of the concept of vague environment in approximate fuzzy reasoning*. *Fuzzy Set Theory and Applications, Tatra Mountains Mathematical Publications, Mathematical Institute Slovak Academy of Sciences, Bratislava, Slovak Republic, vol.12, 1997*, pp. 169-181.
- [13] Kovács, Sz.: *Extending the Fuzzy Rule Interpolation 'FIVE' by Fuzzy Observation*, *Advances in Soft Computing, Computational Intelligence, Theory and Applications*, Bernd Reusch (Ed.), Springer Germany, ISBN 3-540-34780-1, 2006, pp. 485-497. [https://doi.org/10.1007/3-540-34783-6\\_48](https://doi.org/10.1007/3-540-34783-6_48)

- [14] Kovács, Sz.: New Aspects of Interpolative Reasoning. Proceedings of the 6th. International Conference on Information Processing and Management of Uncertainty in Knowledge-Based Systems, Granada, Spain, 1996, pp. 477-482.
- [15] Kovacs, Szilveszter. "Fuzzy Rule Interpolation in Practice." SCIS & ISIS SCIS & ISIS 2006. Japan Society for Fuzzy Theory and Intelligent Informatics, 2006.
- [16] LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. "Deep learning." *nature* 521.7553 (2015): 436-444. <https://doi.org/10.1038/nature14539>
- [17] Li, Yuxi. "Deep reinforcement learning: An overview." arXiv preprint arXiv:1701.07274 (2017). <https://doi.org/10.48550/arXiv.1701.07274>
- [18] Mazyavkina, N., Sviridov, S., Ivanov, S., & Burnaev, E. (2021). Reinforcement learning for combinatorial optimization: A survey. *Computers & Operations Research*, 134, 105400. <https://doi.org/10.1016/j.cor.2021.105400>
- [19] Santra, Santanu, Jun-Wei Hsieh, and Chi-Fang Lin. "Gradient descent effects on differential neural architecture search: A survey." *IEEE Access* 9 (2021): 89602-89618. <https://doi.org/10.1109/access.2021.3090918>
- [20] Tompa, Tamás, and Szilveszter Kovács. "Applying Expert Heuristic as an a Priori Knowledge for FRIQ-Learning." *Acta Polytechnica Hungarica* 17.4 (2020). <https://doi.org/10.12700/aph.17.4.2020.4.2>
- [21] Tompa, Tamás, and Szilveszter Kovács. "Benchmark example for the Heuristically accelerated FRIQ-learning." 2023 24th International Carpathian Control Conference (ICCC). IEEE, 2023. <https://doi.org/10.1109/iccc57093.2023.10178919>
- [22] Tompa, Tamás, and Szilveszter Kovács. "Clustering-based fuzzy knowledgebase reduction in the FRIQ-learning." 2017 IEEE 15th International Symposium on Applied Machine Intelligence and Informatics (SAMI). IEEE, 2017. <https://doi.org/10.1109/sami.2017.7880302>
- [23] Tompa, Tamás, and Szilveszter Kovács. "Determining the minimally allowed rule-distance for the incremental rule-base construction phase of the FRIQ-learning." 2018 19th International Carpathian Control Conference (ICCC). IEEE, 2018. <https://doi.org/10.1109/carpathiancc.2018.8399677>
- [24] Tompa, Tamás, and Szilveszter Kovács. "Heuristically accelerated FRIQ-learning." 20th Jubilee International Symposium on Intelligent Systems and Informatics (SISY 2022). IEEE, 2022. <https://doi.org/10.1109/iccc57093.2023.10178919>
- [25] Tompa, Tamás, and Szilveszter Kovács. "Tudásbázis redukálás a heurisztikusan gyorsított FRIQ-learning rendszerben." *Production Systems and Information Engineering* 11.2 (2023): 1-12. <https://doi.org/10.32968/psaie.2022.4.4>
- [26] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015 <https://doi.org/10.1038/nature14236>
- [27] Vincze, Dávid, Alex Tóth, and Mihoko Niitsuma. "Antecedent redundancy exploitation in fuzzy rule interpolation-based reinforcement learning." 2020 IEEE/ASME International Conference on Advanced Intelligent Mechatronics (AIM). <https://doi.org/10.1109/aim43001.2020.9158875>
- [28] Vincze, Dávid, and Szilveszter Kovács. "Fuzzy rule interpolation-based Q-learning." 2009 5th International Symposium on Applied Computational Intelligence and Informatics. IEEE, 2009. <https://doi.org/10.1109/saci.2009.5136311>
- [29] Vincze, Dávid, and Szilveszter Kovács. "Incremental rule base creation with fuzzy rule interpolation-based Q-learning." *Computational Intelligence in Engineering*. Springer, Berlin, Heidelberg, 2010. 191-203. [https://doi.org/10.1007/978-3-642-15220-7\\_16](https://doi.org/10.1007/978-3-642-15220-7_16)
- [30] Vincze, Dávid, and Szilveszter Kovács. "Rule-base reduction in Fuzzy Rule Interpolation-based Q-learning." *Recent Innovations in Mechatronics* 2.1-2. (2015): 1-6. <https://doi.org/10.17667/riim.2015.1-2/10>