

Genel Ağırlık Atamaları için Verimli Bir Haydut Algoritması

An Efficient Bandit Algorithm for General Weight Assignments

Kaan Gökçesu¹, Tolga Ergen¹, Selami Çiftçi² ve Süleyman S. Kozat¹

¹Elektrik ve Elektronik Mühendisliği Bölümü, İhsan Doğramacı Bilkent Üniversitesi, Ankara, Türkiye

{gokcesu,ergen,kozat}@ee.bilkent.edu.tr

²Turk Telekom Labs, İstanbul, Türkiye

selami.ciftci@turktelekom.com.tr

Özetçe —Bu bildiri, muhalif çok kollu haydut problemi çalışılmış ve genel olarak uygulanabilen verimli bir haydut kolu seçimi yapısı sunulmuştur. Haydut kollarının kayıpları üzerinde hiçbir istatistiksel varsayım yapılmadığı için, bu bildiriye sonuçlar bireysel diziler şeklinde geçerli olmayı garanti etmektedir. Önerilen yapı haydut kolu seçim dizileri üzerinde genel ağırlık atamalarını kullanarak en iyi pişmanlık sınırlarını elde etmektedir. Bu yüzden, bu yapı çok sayıda uygulamada kullanılabilir.

Anahtar Kelimeler—muhalif çok kollu haydut, genel yapı, değiştirme haydutu, verimli uygulama.

Abstract—In this paper, we study the adversarial multi armed bandit problem and present a generally implementable efficient bandit arm selection structure. Since we do not have any statistical assumptions on the bandit arm losses, the results in the paper are guaranteed to hold in an individual sequence manner. The introduced framework is able to achieve the optimal regret bounds by employing general weight assignments on bandit arm selection sequences. Hence, this framework can be used for a wide range of applications.

Keywords—adversarial multi-armed bandit, general framework, switching bandit, efficient implementation.

I. GİRİŞ

Son dönemlerde, çok kollu haydut yaklaşımları birçok gerçek hayat uygulamasında kullanılabildikleri için önemli biçimde çalışılmaktadır [1], [2]. Muhalif çok kollu haydut problemlerinde [3], M tane haydut kolu bulunmakta ve her bir raunt t 'de, olasılıksal olarak bir kol seçilmektedir. Çevrimiçi seçim olan $\{u_t\}_{t \geq 1}$, $u_t \in \{1, 2, \dots, M\}$ verisine dayanarak, sadece seçilen kolun kaybı olan $\{l_{t,u_t}\}_{t \geq 1}$, $l_{t,u_t} \in [0, 1]$ alınmaktadır. Notasyon kolaylığı için $l_{t,u_t} \in [0, 1]$ varsayılmaktadır, ancak, bildiriye çıkarımlar kaydırma ve seviyelendirilmeden sonra herhangi bir sınırlı kayıp için geçerlidir. T raunt bir oyunda, $\mathbf{u}_T = [u_1, \dots, u_T]^T$ sütun vektörü kullanıcının T zamanına kadar yaptığı seçimleri göstermektedir. $\mathbf{s}_T = [s_1, \dots, s_T]^T$ sütun vektörü rastgele olmayan T uzunluğundaki haydut kol seçim dizisini göstermektedir. öyleki her bir t için $s_t \in \{1, 2, \dots, M\}$ geçerlidir. Bildirinin geri

kalan kısmında, \mathbf{s}_T gibi haydut kol seçim dizileri bir strateji olarak belirtilmektedir. $\mathbf{l}_{\mathbf{s}_T} = [l_{1,s_1}, \dots, l_{T,s_T}]^T$ ise \mathbf{s}_T stratejisinin kayıp dizisini göstermektedir. Böylece, \mathbf{u}_T dizisinin kaybı $\mathbf{l}_{\mathbf{u}_T} = [l_{1,u_1}, \dots, l_{T,u_T}]^T$ olmaktadır. Burada, haydut kollarının davranışları üzerinde herhangi bir istatistiksel model varsayılmayan bir muhalif haydut ortamında çalışılmaktadır [4] ve önerilen algoritmalar bireysel dizi şeklinde çalışmayı garanti etmektedir. Algoritmanın t anındaki çıkışı olan u_t kesin bir biçimde çevrimiçi ve rastgeledir, ve çıkış sadece aşağıdaki gibi geçmişteki seçimler ve gözlemlenmiş kayıpların fonksiyonudur:

$$u_t \triangleq u_t(\mathbf{l}_{\mathbf{u}_{t-1}}; \mathbf{u}_{t-1}), \quad u_t \in \{1, \dots, M\}. \quad (1)$$

Herhangi bir \mathbf{s}_T stratejisinin birikimli kayıp fonksiyonu $L_{\mathbf{s}_T} = \sum_{t=1}^T l_{t,s_t}$ olarak gösterilmektedir. Kayıp dizisi üzerinde bir varsayım yapılmadığından, en iyi strateji $\mathbf{s}_T^* = [s_1^*, \dots, s_T^*]^T$ 'ye göre performans aşağıdaki gibi tanımlanmaktadır:

$$\mathbf{s}_T^* = \arg \min_{\mathbf{s}_T} L_{\mathbf{s}_T} \quad \text{veya} \quad s_t^* = \arg \min_{s_t} l_{t,s_t}, \quad 1 \leq t \leq T. \quad (2)$$

Performansı tanımlamak için pişmanlık kavramı şu şekildedir:

$$R_T \triangleq \sum_{t=1}^T l_{t,u_t} - \sum_{t=1}^T l_{t,s_t^*} = L_{\mathbf{u}_T} - L_{\mathbf{s}_T^*} \quad (3)$$

T anına kadar birikmiş pişmanlıktır.

Herhangi bir kayıp dizisi üzerinden en iyi stratejinin performansını elde etmek için M^T stratejinin her biri deterministik uzman olarak düşünülüp, üstel performans ağırlıkları ile birleştirilebilir [5]. Ancak, karışım algoritmaları $O(\sqrt{T \log N})$ pişmanlığı ve $O(N)$ hesaplama karmaşıklığına sahip olduğu için bu üstel sayıda algoritmaların basit bir birleşimi, üstel zamanda yok olmayan bir pişmanlık sınırı $O(T)$ üretmektedir [5]. Bu yüzden, çokterimli zaman içinde kaybolan bir pişmanlığı çevrimiçi olarak elde etmek için stratejilerin akıllıca birleştirilmesi ve ağırlıkların dikkatli ve verimli bir şekilde seçilmesi gerekmektedir. Bunu elde etmek için her bir stratejiye, karmaşıklık maliyetine bağlı olarak değişik bir ağırlık atanmaktadır. Bu ağırlık seçimi AIC ve MDL'nin [6], [7] karmaşıklık cezası ile aynı doğrultudadır.

II. KABA KUVVET YAKLAŞIMI VE PİŞMANLIK SINIRLARI

Bu bölümde, T uzunluğundaki bir oyun için M^T tane olası stratejilerin hepsinin paralel olarak çevrimiçi bir şekilde çalıştırıldığı varsayılmaktadır. Ancak, t anında, M^t tane paralel çalışan strateji bulunmaktadır. Bunlardan her biri farklı bir haydut kolunun kullanılmasını önermektedir. Bu stratejilerden her bir \mathbf{s}_t için o stratejiye güveni gösteren w_{s_t} ağırlığı atanmaktadır. Bu ağırlıklara bağlı olarak, bir ihtimal simpleksi oluşturulmakta ve paralel çalışan her bir stratejiye ağırlıklar normalleştirilerek aşağıdaki gibi bir ihtimal değeri atanmaktadır:

$$P_{s_t} = \frac{w_{s_t}}{\sum_{s'_t \in \mathbb{M}^t} w_{s'_t}}, \quad (4)$$

denkleminde, \mathbb{M}^t t anına kadarki stratejilerin sınıfıdır, ve büyüklüğü $|\mathbb{M}^t| = M^t$ şeklindedir. Her bir kol m için t anında seçim yapabilmek için t anında m 'i öneren stratejiler bulunup onların olasılıkları aşağıdaki gibi toplanmaktadır:

$$p_{t,m} = \sum_{s_t(i:t)=m} P_{s_t}, \quad (5)$$

denkleminde, $\mathbf{s}_t(i:j)$ vektörü, \mathbf{s}_t 'nin i 'den j 'ye kadar olan elemanlarını temsil etmektedir. Aynı haydut kolunu öneren stratejilerin ihtimallerini toplayarak her bir haydut kolunun t anındaki ihtimali oluşturulmaktadır. Denklem (4) ve (5)'teki hesaplamalar direkt olarak w_{s_t} ağırlıklarına bağlıdır. Ağırlık atanması iki bileşene sahiptir. İlk bileşen olarak, her bir \mathbf{s}_t 'ye, sadece s_t 'nin karmaşıklığına bağlı öncül bir ağırlık atanmaktadır. İkinci kısım direkt olarak \mathbf{s}_t 'nin geçmiş performansı olan $\exp(-\eta \tilde{L}_{s_t(1:t-1)})$ üstel ağırlığına bağlıdır. Böylece, birleştirilmiş ağırlıklar şu şekildedir:

$$w_{s_t} = \mathcal{T}(\mathbf{s}_t) e^{-\eta \tilde{L}_{s_t(1:t-1)}}, \quad (6)$$

denkleminde, η öğrenme hızıdır ve $\tilde{L}_{s_t(1:t-1)}$, $L_{s_t(1:t-1)}$ değişkeninin tarafsız tahmincisidir. Burada t anında sadece seçilen kolun kaybı olan l_{t,u_t} gözlemlenmektedir. Bu yüzden, diğer haydut kollarının kaybı için $\tilde{l}_{t,m}$ tahmini olutulmaktadır. Bu amaç için iyi bilinen tarafsız tahminci aşağıdaki gibidir:

$$\tilde{l}_{t,m} = \begin{cases} l_{t,m}/p_{t,m} & m = u_t \\ 0 & m \neq u_t \end{cases} \quad (7)$$

denkleme göre, tahminin beklenen değeri gerçek değer eşittir $\mathbf{E}[\tilde{l}_{t,m}] = l_{t,m}$ [4]. Ayrıca burada, haydut kolları $m \in \{1, \dots, M\}$ üzerinde \mathbf{E}_m tanımlanmaktadır öyleki $\mathbf{E}_m[f(m)] = \sum_{m=1}^M p_{t,m} f(m)$ sağlanmaktadır. Böylece, $\mathbf{E}_m[\tilde{l}_{t,m}] = \sum_{m=1}^M p_{t,m} \tilde{l}_{t,m} = l_{t,u_t}$ eşitliği elde edilmektedir. Birleşim ağırlıkları $\mathcal{T}(\mathbf{s}_t)$ önceden belirlenmektedir. Tam anlamıyla çevrimiçi bir algoritma elde edebilmek için ardışık olarak hesaplanan $\mathcal{T}(\mathbf{s}_t)$ değerleri öyle bir seçilmektedirki iç içe geçme kuralı $\mathcal{T}(\mathbf{s}_t) = \mathcal{T}(\mathbf{s}_t | \mathbf{s}_t(1:t-1)) \mathcal{T}(\mathbf{s}_t(1:t-1))$ elde edilmektedir. Burada, $\mathbf{s}_t(1:t-1)$ stratejisinden \mathbf{s}_t stratejisine ağırlık güncellemesi $\mathcal{T}(\mathbf{s}_t | \mathbf{s}_t(1:t-1))$ ile gösterilmiştir. İhtimal skoru elde edebilmek için göreceli ağırlık güncellemeleri aşağıdaki eşitliği sağlayacak şekilde tasarlanmaktadır:

$$\sum_{m=1}^M \mathcal{T}([\mathbf{s}_t; m] | \mathbf{s}_t) = 1, \quad \forall \mathbf{s}_t, \quad t \in \{0, \dots, T-1\}. \quad (8)$$

Yukarıda, $[\mathbf{s}_t; m]$ vektörü \mathbf{s}_t vektörü ve m ($t+1$ uzunluğunda yeni bir strateji oluşturma) birleştirilmiş halidir, ayrıca,

$\mathbf{s}_0 = [\emptyset]^T$ ve $\mathcal{T}(\mathbf{s}_0) = 1$ geçerlidir. Üstel ağırlıklar her bir stratejinin $t-1$ anına kadarki üstel kayıplarıdır. Bu yüzden, her bir stratejiye atanmış ortak ağırlık w_{s_t} ardışık olarak oluşturulabilir öyleki

$$w_{s_t} = w_{s_t(1:t-1)} \mathcal{T}(\mathbf{s}_t | \mathbf{s}_t(1:t-1)) e^{-\eta \tilde{L}_{s_t(1:t-1)}} \quad (9)$$

sağlanmaktadır.

Denklem (4), (5) ve (6)'daki haydut kol seçim olasılıkları kullanılarak, aşağıdaki pişmanlık sonucu elde edilmektedir.

Teorem 1: $m \in \{1, \dots, M\}$ 'in muhalif çok kollu haydutun kolları olduğunu ve $l_{t,m} \in [0, 1]$ kaybının t anında m kolunu seçmekten dolayı oluşan kayıp olduğunu varsayalım. Denklem (8)'yi sağlayan ardışık olarak oluşturabilen herhangi bir birleşim ağırlık ataması $\mathcal{T}(\cdot)$ ve her bir kolun seçin ihtimalleri belirlemek için olan (6)'daki gibi üstel kayıpları kullanılarak, aşağıdaki beklenen pişmanlık elde edilmektedir:

$$\mathbf{E}[R_T] \leq \min_{s_T} \left(\frac{\eta M T}{2} + \frac{1}{\eta} \ln W(s_T) + L_{s_T} - L_{s_T^*} \right) \quad (10)$$

T raunt oyunda el edilmiştir. Burada, $\eta \geq 0$ üstel ağırlıklar-daki öğrenme hızını ve $W(s_T) \triangleq 1/\mathcal{T}(s_T)$, s_T stratejisinin birleşim ağırlıklarının tersidir. s_T stratejisinin birimli kaybı L_{s_T} ve s_T^* en iyi kol seçim stratejisinin birimli kaybı $L_{s_T^*}$ ile gösterilmektedir. En iyi strateji bütün kolların, $m \in \{1, \dots, M\}$, bütün zaman indekslerindeki, $t \in \{1, \dots, T\}$, kayıpları bilindiği öncül bilgisi ile seçilmiştir.

Teorem 1'deki sonuç dikkatli bir şekilde tasarlanmış $\mathcal{T}(\cdot)$ ve η ile altdoğrusal ve hatta en uygun pişmanlığın elde edilebileceğini göstermektedir. Ancak, ağırlık ataması $\mathcal{T}(\cdot)$ 'nin ardışık olarak oluşturabilir olması ve (8)'i sağlaması gerekmektedir. Buna ek olarak, Teorem 1'deki sonuç önerilen yapının performansının en iyi stratejinin karmaşıklık maliyetinin ($W(s_T^*)$) yanısıra kaybı en iyi stratejinin kaybına (en iyi kayıp) göreceli olarak yakın olan stratejilerin karmaşıklık maliyetine de bağlı olduğunu göstermektedir. Bu yüzden, en iyi strateji yüksek karmaşıklık maliyetine sahip olsa bile, eğer en iyi kayba yeterince yakın düşük karmaşıklık maliyeti olan bir strateji varsa, önerilen algoritma göreceli olarak düşük bir pişmanlık elde edebilir. Denklem (10)'daki beklenti, sonuçların istatistiksel varsayımı olmayan haydut kayıplarının herhangi bir dizisi için eşit oranda geçerli olması için sahip olunan rastgelelik yüzündendir.

Teorem 1'in ispatı: t anında, en iyi seçim stratejisi \mathbf{s}_T^* 'ye karşı olan pişmanlık $r_t = l_{t,u_t} - l_{t,s_t^*}$ ile gösterilmektedir. r_t daha idare edilebilir bir forma sokulup, iki farklı terim elde edilmektedir. Bu terimler ayrı ayrı olarak aşağıdaki gibi sınırlandırılmaktadır:

$$r_t = \left(l_{t,u_t} + \frac{\ln \mathbf{E}_m[e^{-\eta \tilde{l}_{t,m}}]}{\eta} \right) - \left(\frac{\ln \mathbf{E}_m[e^{-\eta \tilde{l}_{t,m}}]}{\eta} + l_{t,s_t^*} \right). \quad (11)$$

Denklem (11)'deki ilk terim $x > 0$ için $\ln x \leq x-1$ kullanılarak aşağıdaki gibi sınırlandırılmaktadır:

$$\ln \mathbf{E}_m[e^{-\eta \tilde{l}_{t,m}}] \leq \mathbf{E}_m[e^{-\eta \tilde{l}_{t,m}} - 1]. \quad (12)$$

Denklem (12), $x > 0$ için $e^{-x} - 1 + x \leq x^2/2$ kullanılarak aşağıdaki ifade elde edilmektedir:

$$\ln \mathbf{E}_m[e^{-\eta \tilde{l}_{t,m}}] \leq \mathbf{E}_m \left[\frac{\eta^2 \tilde{l}_{t,m}^2}{2} \right] - \eta \mathbf{E}_m[\tilde{l}_{t,m}] \leq \frac{\eta^2 l_{t,u_t}^2}{2 p_{t,u_t}} - \eta l_{t,u_t}. \quad (13)$$

Denklem (13) denklem (11)'deki ilk terimde yerine konulursa

$$r_t \leq \frac{\eta}{2p_{t,u_t}} + \left[-\frac{1}{\eta} \ln \mathbf{E}_m [e^{-\eta \tilde{L}_{t,m}}] - l_{t,s_t^*} \right], \quad (14)$$

$l_{t,u_t} \leq 1$ olduğu için üstteki denklem elde edilmektedir. Denklem (14)'teki ikinci terimi üstten sınırlamak için (5) ve (4) kullanılarak aşağıdaki gibi beklenti hesaplanmaktadır:

$$\begin{aligned} \mathbf{E}_m [e^{-\eta \tilde{L}_{t,m}}] &= \sum_{m=1}^M p_{t,m} e^{-\eta \tilde{L}_{t,m}} = \sum_{s_t^* \in \mathbb{M}^t} P_{s_t^*} e^{-\eta \tilde{L}_{t,s_t^*(t:t)}}, \\ &= \sum_{s_t^* \in \mathbb{M}^t} \frac{w_{s_t^*}}{\sum_{s_t' \in \mathbb{M}^t} w_{s_t'}} e^{-\eta \tilde{L}_{t,s_t^*(t:t)}}, \\ &= \frac{\sum_{s_t^* \in \mathbb{M}^t} \mathcal{T}(s_t^*) e^{-\eta \tilde{L}_{s_t^*}}}{\sum_{s_t' \in \mathbb{M}^t} \mathcal{T}(s_t') e^{-\eta \tilde{L}_{s_t'}}}, \quad (15) \end{aligned}$$

burada, (6) ve (8) kullanılmaktadır. Bundan sonra, (15)'in logaritmasını bütün T rauntları için toplayıp, $-1/\eta$ ile çarparak aşağıdaki denklem elde edilmektedir:

$$\begin{aligned} \sum_{t=1}^T -\frac{1}{\eta} \ln \mathbf{E}_m [e^{-\eta \tilde{L}_{t,m}}] &= -\frac{1}{\eta} \ln \sum_{s_t^* \in \mathbb{M}^T} \mathcal{T}(s_t^*) e^{-\eta \tilde{L}_{s_t^*}}, \\ &\leq -\frac{1}{\eta} \ln [\mathcal{T}(s_T) e^{-\eta \tilde{L}_{s_T}}] \leq -\frac{1}{\eta} \ln \mathcal{T}(s_T) + \tilde{L}_{s_T}. \quad (16) \end{aligned}$$

Yukarıdaki denklem herhangi bir $s_T \in \mathbb{M}^T$ için geçerlidir. Denklem (14) bütün T rauntları için toplanarak aşağıdaki gibi T rauntluk oyunda birikmiş pişmanlık değeri bulunmaktadır:

$$R_T = \sum_{t=1}^T \frac{\eta}{2p_{t,u_t}} - \frac{1}{\eta} \sum_{t=1}^T \ln \mathbf{E} [e^{-\eta \tilde{L}_{t,m}}] - \sum_{t=1}^T l_{t,s_t^*}. \quad (17)$$

Denklem (16)'yı denklem (17)'de kullanarak, toplam pişmanlık

$$R_T \leq \sum_{t=1}^T \frac{\eta}{2p_{t,u_t}} - \frac{1}{\eta} \ln \mathcal{T}(s_T) + \tilde{L}_{s_T} - L_{s_T^*} \quad (18)$$

şeklinde yazılır. Seçim u_t , ve böylece p_{t,u_t} , R_T 'deki rastgele değişkenlerdir. Denklem (18)'in kol seçim ihtimallerine göre beklentisi alındığında

$$\mathbf{E}[R_T] \leq \frac{\eta MT}{2} - \frac{1}{\eta} \ln \mathcal{T}(s_T) + L_{s_T} - L_{s_T^*}$$

elde edilmektedir. Gösterim kolaylığı için birleşim ağırlıklarından stratejinin karmaşıklık maliyetine kadar olan gösterimler $W(s_T) \triangleq 1/\mathcal{T}(s_T)$ olacak şekilde değiştirilmektedir. Böylece,

$$\mathbf{E}[R_T] \leq \frac{\eta MT}{2} + \frac{1}{\eta} \ln W(s_T) + L_{s_T} - L_{s_T^*} \quad (19)$$

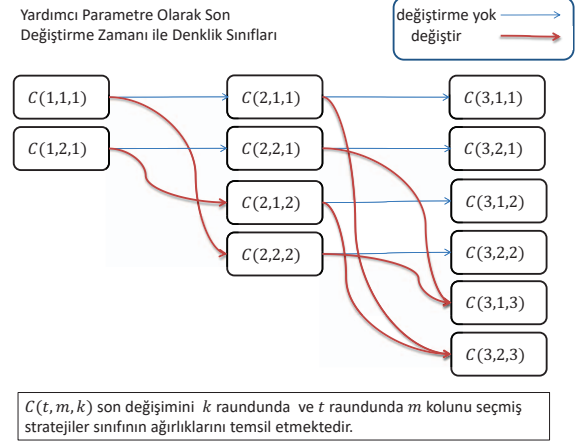
elde edilmektedir. Denklem (19) herhangi bir strateji s_T için sağlandığı için daha sıkı bir sınır (19)'u s_T üzerinden en küçük duruma getirerek elde edilebilir. Ve bu işlem (10)'u verir. \square

Sonuç 1: Kayıp dizisinden bağımsız bir üst sınır elde etmek için Teorem 1'de $s_T = s_T^*$ eşitliği kullanılarak (10) denklemi aşağıdaki gibi üstten sınırlandırılmaktadır:

$$\mathbf{E}[R_T] \leq \frac{\eta MT}{2} + \frac{1}{\eta} \ln W(s_T^*). \quad (20)$$

Yukarıdaki denklemde, $W(s_T^*)$ en iyi kol seçim stratejisi s_T^* 'nin birleşim ağırlıklarının tersidir.

III. VERİMLİ UYGULAMA



Şekil 1: Son değiştirme zamanını yardımcı parametre olarak kullanarak, iki kolu haydut durumunun ilk 3 raundu için verimli birleşim örneği. Bu durumda, yardımcı parametre vektörü σ_t sadece k 'yi içermektedir ve alabileceği olası değerler zaman ile doğrusal olarak artmaktadır. Bu yüzden, k sadece t değişik değer alabildiği için bu yapı doğrusal karmaşıklığa sahip bir algoritmayı formüllestirmektedir.

Hesaplama karmaşıklığını azaltmak için belli stratejileri beraber gruplayarak denklik sınıfları oluşturulmaktadır. $C(t, m, \sigma_t)$ m kolunun denklik sınıfının ve t anındaki σ_t yardımcı parametresinin ağırlığı olarak tanımlanmaktadır. Denklik sınıfı $C(t, m, \sigma_t)$ t anında m kolunu seçen bütün s_t stratejilerini içermektedir ve davranışı σ_t parametre vektörü ile eşleşmektedir. Örnek olarak, Şekil 1'deki σ_t vektörünün sadece stratejilerin yaptığı son değiştirmenin zaman indeksini içermesi düşünülebilir. Stratejileri son değiştirme zamanına göre gruplandırmak, doğrusal olarak artan sayıda denklik sınıflarına neden olmaktadır. Yardımcı parametre σ_t farklı şekilde gruplar da içerebilir. Örnek olarak, stratejilerin yaptığı değişim sayısı verilebilir.

σ_t 'ya dahil edilen parametreler, onun en sonda kaç tane strateji belirleyeceğini ve kaç tane denklik sınıfına sahip olacağını belirlemektedir. Yardımcı parametre σ_t kullanılmasının sebebi (9)'daki ağırlık güncellemeleri aynı olan belli stratejileri gruplamaktır. Bu yüzden, σ_t 'ya birleşim ağırlık güncellemeleri, $\mathcal{T}(s_t | s_t(1:t-1))$, ile alakalı bütün parametreler dahil edilmelidir. Böylece, bileşim ağırlık atanmasının, $\mathcal{T}(\cdot)$, tasarlanması σ_t 'ya dahil edilecek parametreleri etkilemektedir. Burada, Λ_t olası bütün σ_t vektörlerini içeren vektör uzayı olarak tanımlanmaktadır.

Denklik sınıfının ağırlığı, sınıf parametreleri t, m, σ_t ile uyumlu davranan stratejilerin ağırlıkları toplamıdır. Buna göre aşağıdaki denklem elde edilmektedir.

$$C(t, m, \sigma_t) = \sum_{\substack{s_t(t:t)=m \\ \sigma(s_t)=\sigma_t}} w_{s_t}. \quad (21)$$

Yukarıda, $\sigma(\cdot)$ s_t 'den σ_t parametresine kadar olan eşleme fonksiyonudur ve $\sigma : \mathbb{M}^t \rightarrow \Lambda_t$ olarak tanımlanmaktadır. Ayrıca, w_{s_t} (6)'da tanımlanmaktadır. Şekil 1'de, iki kollu haydut oyununun ilk üç raunda için denklik sınıfı örneği verilmiştir. Bu şekilde, $C(t, m, k)$ son değiştirmeyi k raundunda yapmış ve t raundunda m kolunu seçmiş stratejilerin ağırlıklarının temsil etmektedir. Örnek olarak, $C(3, 1, 3)$ $s_3 \in \{[2, 2, 1]^T, [1, 2, 1]^T\}$ stratejiler sınıfının ağırlığıdır. Buradaki amaç (9)'daki çarpımsal güncellemeyi belli sayıda stratejiler için aynı anda yapmak olduğu için, $\mathcal{T}(s_t | s_t(1:t-1)) \exp(-\eta \tilde{l}_{t,s_t(1:t-1)})$ güncellemesinin karışımındaki bütün stratejiler için aynı olması gerekmektedir. Üstel kayıp güncellemesi eğer stratejiler tarafından seçilen kollar aynı ise aynıdır. Bu koşul ancak stratejiler aynı denklik sınıfına ait olduğunda gerçekleşmektedir. Birleşim ağırlıkları güncellemeleri sınıf parametreleri olan şimdiki raunt t , şimdiki raunttaki kol seçimi m ve yardımcı parametre σ_t 'ya bağlı olarak tasarlanmaktadır öyleki aynı sınıfta olan stratejiler aynı ağırlık güncellemesine sahiptirler. Denklik sınıfı $C(t, m', \sigma_t)$ 'den $C(t+1, m, \sigma_{t+1})$ 'ye kadar olan ortak birleşim ağırlık güncellemesi $\mathcal{T}(t+1, m, \sigma_{t+1} | t, m', \sigma_t)$ ile gösterilmektedir. Bu gösterimde, m' ve σ_t sonraki zaman indeksleri arasında ayırım yapmak için kullanılmaktadır. Bu yüzden,

$$C(t+1, m, \sigma_{t+1}) = \sum_{m', \sigma_t} C(t, m', \sigma_t) \mathcal{T}(t+1, m, \sigma_{t+1} | t, m', \sigma_t) e^{-\eta \tilde{l}_{t, m'}} \quad (22)$$

her bir denklik sınıfı ağırlığı kendi parametrelerine uyan stratejilerin ortak ağırlığının toplamı olarak hesaplandığı için yukarıdaki denklem geçerli olmaktadır.

Üstel kayıp güncellemesi son seçilen kola bağlı olduğu için, denklik sınıfları stratejileri son seçtikleri kola göre gruplamaktadır. σ_t 'daki yardımcı parametre birleşim ağırlıklarını güncellemek için kullanılmaktadır. Algoritma 1'de, genel yapının tam verimli uygulaması sunulmaktadır.

Denklem (22), (21)'i kullanılarak (9) denkleminin direkt bir uygulaması olduğu için, Algoritma 1'deki verimli uygulama, direkt olarak Bölüm II'deki kaba kuvvet yaklaşımının ağırlık atamasını uygulamaktadır. Bu yüzden, kaba kuvvet yaklaşımı için yapılan bütün pişmanlık analizleri (Teorem 1 ve Sonuç 1 gibi) Algoritma 1'deki verimli uygulama için de geçerlidir. Hesaplama karmaşıklığı denklik sınıfı sayısı (t anında $M|\Lambda_t|$ ile gösterilmektedir.) ile alakalı olduğu için, denklik sınıfları kullanılarak hesaplama karmaşıklığı zaman içinde üstel olmaktan terimsel olmaya indirgenebilmektedir.

Açıklama 1: Yardımcı parametre σ_t kullanılarak genel ağırlık atamasında daha fazla esneklik sağlanmaktadır. Bu yapının genelliği ve ağırlık ataması çeşitli uygulamalar için bir çok ihtimal sağlamaktadır. Değişik ortamlar için değişik ağırlık atamaları tasarlanabilir. Ağırlık atamaları değişik karmaşıklıkta maliyet fonksiyonlarına uygun hale getirilebilir. Örnek olarak, bütün değişimlere eşit davranmak yerine uzun bölümlerden sonraki değişimlere daha fazla önem verilmesi düşünülebilir. Eğer bölümler belli uzunluktan kısa ise aykırı olarak düşünülüp, değişim olarak görülmeyebilir. Bu örnek için son değişim zamanı yardımcı değişken olarak kullanılması ile uygun bir ağırlıklandırma şeması tasarlanabilir. Buna ek olarak, bu genel yapı, bütün küme \mathbb{M}^T yerine sadece stratejilerin makul bir alt kümesini birleştirmek için kullanılabilir. Örnek olarak, eğer en iyi kol en azından K raunt için değişmiyorsa, son parçanın uzunluğu yardımcı değişken olarak kullanılabilir ve parçalar

Algorithm 1 Verimli Genel Yapı

```

1: Sabit  $\eta \in \mathbb{R}^+$ 'ya ilk değer ata
2: Birleşim ağırlık atamalarını seç
3:  $t \in 1, \dots, T$  için  $\Lambda_t$ 'yi belirle
4:  $\Lambda_1$ 'in uzantısı olan  $\sigma_1$ 'ya ilk değer atama
5:  $m \in 1, \dots, M$  için  $C(1, m, \sigma_1) = 1/M$ 'yi belirle
6:  $p_{1,m} = C(1, m, \sigma_1)$  için ilk değer ata
7: for  $t = 1 : T$  do
8:    $p_{t,m}$  ihtimali ile  $M$  koldan birini seç
9:   Kayıp  $l_{t,u_t}$ 'nin alınması
10:   $m \in 1, \dots, M$  için  $\tilde{l}_{t,m} = \frac{l_{t,m} \mathbb{1}_{m=u_t}}{p_{t,m}}$ 'yi belirle
11:  for  $\sigma_{t+1} \in \Lambda_{t+1}$  do
12:    for  $m = 1 : M$  do
13:       $C(t+1, m, \sigma_{t+1}) = \sum_{m', \sigma_t} C(t, m', \sigma_t) \mathcal{T}(t+1, m, \sigma_{t+1} | t, m', \sigma_t) e^{-\eta \tilde{l}_{t, m'}}$ 
14:    end for
15:  end for
16:  for  $m = 1 : M$  do
17:     $p_{t+1, m} = \frac{\sum_{\sigma_{t+1} \in \Lambda_{t+1}} C(t+1, m, \sigma_{t+1})}{\sum_{m=1}^M \sum_{\sigma_{t+1} \in \Lambda_{t+1}} C(t+1, m, \sigma_{t+1})}$ 'yi belirle
18:  end for
19: end for

```

K uzunluğuna ulaşmadan olan değişimleri engelleyerek, sadece en azından K uzunluğundaki parçaya sahip stratejiler birleştirilebilir. Eğer belli bir kol m en iyi kol olan m' 'dan hemen sonraki en iyi kol olmazsa, sadece makul stratejileri birleştirmek için m' 'dan m 'e olan değişimleri engelleyen bir ağırlıklandırma şeması tasarlanabilir.

IV. SONUÇLAR

Bu bildiriye, muhalif çok kollu haydut problemi çalışılmış ve genel olarak uygulanabilen verimli bir haydut kol seçim yapısı önerilmiştir. Çeşitli uygulamalar için, önerilen yapı her türlü ağırlıklandırma şeması ile çalışmaktadır. Bu yapı ardışık olarak bütün olası kol seçim stratejilerini dikkatli bir şekilde oluşturulmuş ağırlıklar ile birleştirmektedir. Burada, olası strateji sayısı M^T ile büyümektedir. Belirli stratejileri gruplandırıp topluca güncelleyen denklik sınıfları yaratarak, bu ağ yapısı verimli bir şekilde uygulanmıştır.

KAYNAKLAR

- [1] V. Krishnamurthy and R. J. Evans, "Hidden markov model multiarm bandits: a methodology for beam scheduling in multitarget tracking," *Signal Processing, IEEE Transactions on*, vol. 49, no. 12, pp. 2893–2908, 2001.
- [2] —, "Correction to "hidden markov model multiarm bandits: a methodology for beam scheduling in multitarget tracking"," *IEEE Transactions on Signal Processing*, vol. 51, no. 6, pp. 1662–1663, 2003.
- [3] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire, "Gambling in a rigged casino: The adversarial multi-armed bandit problem," in *Foundations of Computer Science, 1995. Proceedings., 36th Annual Symposium on*, Oct 1995, pp. 322–331.
- [4] S. Bubeck and N. Cesa-Bianchi, "Regret analysis of stochastic and nonstochastic multi-armed bandit problems," *Foundations and Trends in Machine Learning*, vol. 5, no. 1, pp. 1–122, 2012.
- [5] P. Auer, N. Cesa-Bianchi, Y. Freund, and R. E. Schapire, "The nonstochastic multiarmed bandit problem," *SIAM J. Comput.*, vol. 32, no. 1, pp. 48–77, Jan. 2003.
- [6] H. Akaike, "A new look at the statistical model identification," *IEEE Transactions on Automatic Control*, vol. 19, no. 6, pp. 716–723, Dec 1974.
- [7] J. Rissanen, "Modeling by shortest data description," *Automatica*, vol. 14, no. 5, pp. 465–471, 1978.