

RASSAL ÇİZGE SERİLERİNDE ÇOKLU DEĞİŞİM NOKTASI ANALİZİ MULTIPLE CHANGE POINT ANALYSIS IN RANDOM GRAPH SERIES

Türkan Hamzaoğlu¹, Barış Kurt¹, A. Taylan Cemgil¹

¹Bilgisayar Mühendisliği Bölümü
Boğaziçi Üniversitesi, İstanbul

turkan.hamzaoglu@boun.edu.tr, baris.kurt@boun.edu.tr, taylan.cemgil@boun.edu.tr

ÖZETÇE

Çizgeler süreçlerin modellenmesinde kullanılan önemli matematiksel araçlardır. Bu alanda önemli bir problem, zamanla üretici süreç parametrelerinde meydana gelen değişimleri çıkarımlamaktır. Bu çalışmada, rassal öbek çizgeleri serilerinde çoklu değişim noktalarını çıkarımlama problemi ele alındı. Zaman serilerinde çıkarım yapmak için en sık kullanılan algoritmalarından biri ileri-geri algoritmasıdır. Bu algoritmanın çizge değerli modellerde hesaplama karmaşıklığını azaltmak için, geri yönlü mesaj kısmı değiştirilip geriye doğru Monte Carlo örnekleme uyarlandı. Yapılan testlerle tasarlanan algoritmanın gerçeğe uygun çıkarımlar yaptığı gözlemlendi.

ABSTRACT

Graphs are important mathematical tools for modelling processes. An important issue in this area is to infer the changes that occur in the underlying generative process. In this work, inference of multiple change points in stochastic block graph time series is studied. A well-known algorithm for inference in time series is the forward-backward algorithm. In order to decrease computational complexity of this algorithm in graphical models, backward smoothing part is replaced with backward Monte Carlo sampling. With the experiments, it is observed that modified algorithm gives result in accordance with the real data.

1. GİRİŞ

Rassal çizgeler, birçok alanda process ve veri modellemelerinde ve veri analizinde kullanılmaktadır. En yaygın kullanım alanlarının başında telekomünikasyon ağlarının incelenmesi, ulaşım şebekesinin modellenmesi, mobil telefon ağlarının incelenmesi, proteinler arası etkileşimlerin modellenmesi, besin zinciri gibi biyolojik ağların analizi; facebook, twitter sosyal ağların analizi; netflix, amazon gibi kullanıcı tercihlerinin incelenmesidir. Bu örnekler daha da çoğaltılabilir.

Yukarıdaki ve benzeri analizler için değişik çizge modelleri kullanılmaktadır. Literatürde en çok yer alan modeller: Erdős-Reyni-Gilbert modeli, sosyal ağlar için P_1 ve P_2 modelleri, üssel rassal çizge modelleri, rassal blok modelleri, tercihli bağlanma modeli, çoğaltmalı bağlanma modeli ve türevleri olarak sayılabilir[1].

Bu makalede, Rassal Blok Model Çizgeleri ve bunların dinamik yani zaman içerisinde değişimi üzerinde çalışıldı. Bu modelde düğümlerin örtülü kategorilere dahil oldukları ve düğümler arasındaki ilişkilerin, bu örtülü kategorilere göre oluştuğu varsayılır[2]. Örnek olarak, bir basit kullanıcı-servis çizgesi düşünelim. Kullanıcıların tercih ettikleri servisler kullanıcı ve servis kategorilerine göre belirlenmektedir. Modelle ilgili matematiksel ayrıntılar ileriki bölümlerde verilecektir.

Günümüzde durağan veri analizi kadar, dinamik veri analizi de önem kazandı. Örnek olarak, ilgili süreçlerdeki değişimi takip etmek ve değişime göre önlem almak ya da eylemde bulunmak ilgili kurumlar için önemlidir. Rassal blok modellerinden örnek vermek gerekirse, zaman içinde müşteri ve servis tercihleri arasındaki değişimi fark etmek, servis sağlayıcıları için önemli olacaktır. Bu ve benzeri zaman içindeki değişim noktalarını çıkarımlamak için rassal çizge modelleri zaman serisi modelleri olarak tasarlandı.

2. RASSAL ÖBEK ÇİZGELERİ ve ÇOKLU DEĞİŞİM NOKTALARI

Bu bölümde önce rassal öbek çizgelerinin matematiksel modeli incelenecek, sonrasında da rassal öbek çizgeleri zaman serilerinde çoklu değişim noktalarının nasıl modellendiğine değineceğiz. Çoklu değişim noktalarının çıkarımlanmasında ise saklı Markov modellerinin analizinde kullanılan *ileri yönde filtreleme ve geri yönde düzeltme* (forward filtering-backward smoothing) algoritmasının Monte Carlo örneklemesine uyarlanması olan *ileri yönde filtreleme ve geri yönde örnekleme* (forward filtering-backward sampling) algoritması kullanılmıştır.

Bu çalışma Boğaziçi Üniversitesi BAP Koordinatörlüğü 5723 numaralı proje tarafından desteklenmektedir.

2.1. Rassal Öbek Çizgeleri

Temel rassal öbek çizgelerinde herbir düğüm örtülü alt bir kategoriye aittir. Düğümler arası bağlantılar da bu kategorilere göre belirlenir. Literatürde, bu modelin çeşitlemeleri tanımlanmıştır. Bazı modellerde, düğümler arasındaki kenarlar alt kategorilere ek olarak, düğümlere ya da düğümlerin çeşitli özelliklerine de bağlanabilir. Bu makalede, temel rassal öbek çizgelerini ele alacağız ama yukarıda belirtilen çeşitlemeler de kolayca uyarlanabilir. Şekil 1’de modelin Bayeşçi çizgelerle betimlenmesi gösterilmektedir. Önceki kullanıcı-servis örneği baz alınarak, altta yatan matematiksel model şu şekilde betimlenebilir:

$N^u = \text{Kullanıcı sayısı}$

$K^u = \text{Kullanıcı kategori sayısı}$

$N^s = \text{Servis sayısı}$

$K^s = \text{Servis kategori sayısı}$

$k \in \{1, \dots, K^u\} = \text{Kullanıcı kategori endeksi}$

$l \in \{1, \dots, K^s\} = \text{Servis kategori endeksi}$

$p \in \{1, \dots, N^u\} = \text{Kullanıcı endeksi}$

$q \in \{1, \dots, N^s\} = \text{Servis endeksi}$

$C_p^u \sim \text{Multn}(\pi^u) = p \text{ kullanıcısının kategorisi}$

$C_q^s \sim \text{Multn}(\pi^s) = q \text{ servisinin kategorisi}$

$B_{kl} \sim \text{Beta}(\alpha, \beta) = k \text{ kategorili kullanıcı ile } l \text{ kategorili servis arasında ilişki olasılığı}$

$$Y_{pq} = \begin{cases} 1 & \text{eğer } p \text{ kullanıcısı } q \text{ servisini kullanıyorsa} \\ 0 & \text{diğer şartlarda} \end{cases}$$

$$\sim \text{Bern}(C_p^u B(C_q^s)^T) = \text{kullanıcılar ve servisler arasındaki komşuluk matrisi}$$

(1)

Tüm modelin bileşik olasılık dağılımı şu şekilde elde edilir.

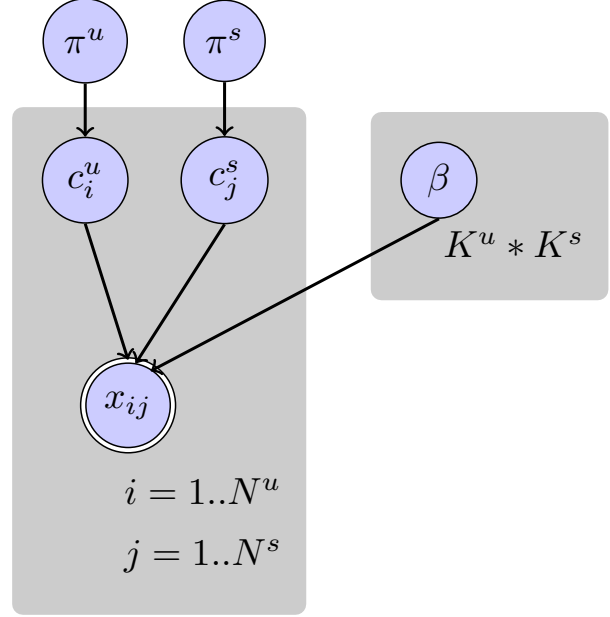
$$p(Y, B, C_{1:N^u}^u, C_{1:N^s}^s) = p(C_{1:N^u}^u) * p(C_{1:N^s}^s) * p(B) * p(Y|\beta, C_{1:N^u}^u, C_{1:N^s}^s)$$

$$= \prod_p^{N^u} \prod_k^{K^u} (\pi_k^u)^{C_{pk}^u} * \prod_q^{N^s} \prod_l^{K^s} (\pi_l^s)^{C_{ql}^s} * \prod_k^{K^u} \prod_l^{K^s} \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} B_{kl}^{(\alpha-1)} (1 - B_{kl})^{(\beta-1)}$$

$$* \prod_p^{N^u} \prod_q^{N^s} \prod_k^{K^u} \prod_l^{K^s} (B_{kl}^{Y_{pq}} (1 - B_{kl})^{(1-Y_{pq})})^{C_{pk}^u * C_{ql}^s}$$

(2)

Bu model için düğümlere kategori atama (C_u ve C_s) ve bağlantı parametrelerini (B) çıkarım yöntemleri geçen yılki makalemizde incelenmişti [3]. Burada dikkat edilmesi gereken önemli nokta, düğümlerin kategorileri bilindiğinde, rassal öbek modeli koşullu bağımsız olarak



Şekil 1: Rassal Öbek Çizge Modeli

$C_u * C_s$ tane alt Erdos-Reyni çizgesi olarak ifade edilebilir.

Herhangi bir A Erdos-Reyni matrisini alalım. Bu matris kenarlar arası bağlantı olasılığı $Bernoulli(b)$ olasılık dağılımı ile belirlenmiş olsun. b parametresinin eşlenik önsel dağılımı da $Beta(\alpha, \beta)$ olsun. Bu durumda A matrisinin marjinal olasılığı aşağıdaki şekilde kolayca bulunabilir:

$$p(A) = \int_b p(b) p(A|b) db \quad (3)$$

$$c = \sum_{i,j} A_{ij} \quad (\text{toplam bağ sayısı}) \quad (4)$$

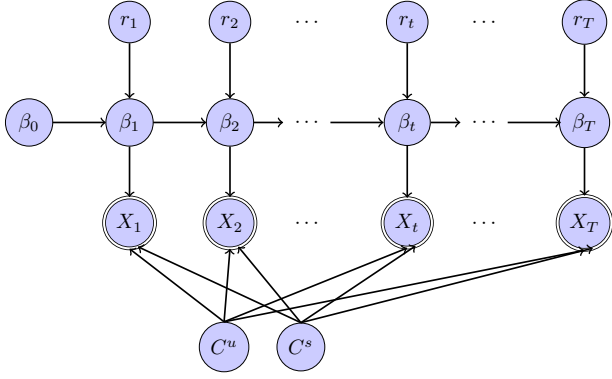
$$n = \sum_{i,j} 1 \quad (\text{olabilecek tüm bağların sayısı}) \quad (5)$$

$$p(b) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} b^{\alpha-1} (1 - b)^{\beta-1} \quad (6)$$

$$\log p(A) = \log \int_0^1 \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} b^{\alpha+c-1} (1 - b)^{\beta+n-c-1} db \quad (7)$$

$$= \log \Gamma(\alpha + \beta) + \log \Gamma(\alpha + c) + \log \Gamma(\beta + n - c) - \log \Gamma(\alpha) - \log \Gamma(\beta) - \log \Gamma(\alpha + \beta + n) \quad (8)$$

Yukarıdaki denklemlere göre, A matrisinin bağlantı parametresi b 'nin çıkarımı oldukça basittir.



Şekil 2: Çoklu Değişim Noktaları Modeli

2.2. Rassal Öbek Çizgilerinin Zaman Serilerinde Çoklu Değişim Noktalarının Çıkarımı

Çoklu değişim noktaları modelinde, T uzunluğundaki zaman serisi için değişim noktalarının sayısını ve yerleri bilinmiyordur ve amaç bu değişim noktalarının yerlerini çıkarılmaktır. Bu modelde, herhangi bir t anında değişim olduysa, rassal öbek çizge modelimizdeki bağlantı parametrelerinin önsel dağılımından tekrar çekildiği, eğer değişim olmadysa $t - 1$ bağlantı parametreleriyle aynı olduğu varsayılır. Varsayılan diğer önemli bir nokta da, öbek çizgelerindeki düğümlerin kategorilerinin bilinmiyor ve T süresince değişmiyor olmasıdır. Her t anında öbek çizgeleri ilgili bağlantı parametrelerine ve altta yatan kategori eşlemesine göre yeniden oluşturuluyor. Şekil 2'de çoklu değişim noktaları modelimizin Bayesçi çizgelerle temsilini görüyoruz.

Matematiksel olarak üretici modeli şu şekilde tanımlayabiliriz:

$$r_t \sim Be(p) \quad (9)$$

$$\beta_t \sim [r_t = 0]\delta(B_t = B_{t-1}) \quad (10)$$

$$B_t \sim [r_t = 0]\delta(B_t = B_{t-1}) + [r_t = 1] \prod_k^{K^u} \prod_l^{K^s} Beta(\alpha, b) \quad (11)$$

$$C_p^u \sim Multn(\pi^u) \quad (12)$$

$$C_q^s \sim Multn(\pi^s) \quad (13)$$

Herhangi bir t anındaki komşuluk matrisinin (Y_t) koşullu olasılığı önceki bölümde verildiği gibi hesaplanır:

$$Y_t | B_t, C^u, C^s = \prod_p^{N^u} \prod_q^{N^s} \prod_k^{K^u} \prod_l^{K^s} (B_{t,kl}^{Y_t,pq}) (1 - B_{t,kl})^{(1-Y_t,pq)} C_{pk}^u * C_{ql}^s \quad (14)$$

2.2.1. İleri-Geri Algoritması

İleri-geri algoritması zaman serileri ve saklı Markov modellerinde anlık parametre çıkarımlarında çok sık kullanılan bir tam çıkarım yöntemidir[4]. Bizim modelimizde de $p(r_t, B_t | Y_{1:T})$ koşullu dağılımları bu algoritma ile hesaplanabilir. Algoritmanın adımlarını kısaca tanımlayalım.

Alfa (ileri yönlü) mesajları:

$$\begin{aligned} \alpha_{0|0} &= p(B_0) \\ t &= 1..T \\ \alpha_{t|t-1} &= p(r_t, B_t, Y_{1:t-1}) \\ \alpha_{t|t} &= p(r_t, B_t, Y_{1:t}) \end{aligned} \quad (15)$$

Beta (geri yönlü) mesajları:

$$\begin{aligned} \beta_{T|T} &= p(Y_T | r_T, B_T) \\ t &= T - 1, ..1 \\ \beta_{t|t+1} &= p(Y_{t+1:T} | r_t, B_t) \\ \beta_{t|t} &= p(Y_{t:T} | r_t, B_t) \end{aligned} \quad (16)$$

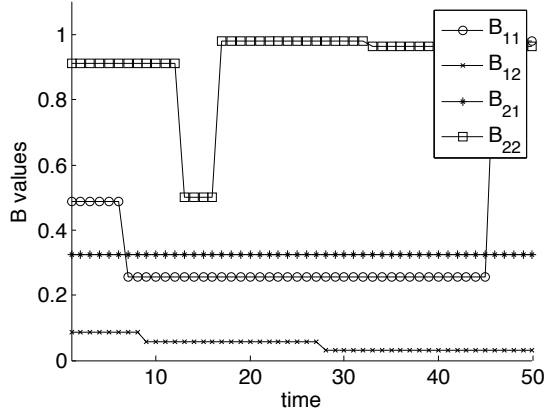
Değişim noktalarının sonsal dağılımları da alfa ve beta mesajları kullanılarak şu şekilde elde edilir:

$$\begin{aligned} p(r_t, B_t | Y_{1:T}) &\sim p(Y_{1:T}, r_T, B_T) \\ &= p(Y_{1:t-1}, r_t, B_t) p(Y_{t:T} | r_t, B_t, Y_{1:t-1}) \\ &= p(Y_{1:t-1}, r_t, B_t) p(Y_{t:T} | r_t, B_t) \\ &= \alpha_{t|t-1} \beta_{t|t} \end{aligned} \quad (17)$$

2.2.2. İleri Yönlü Filtreleme - Geri Yönlü Örnekleme Algoritması

İleri-geri algoritması çıkarılmak istenen parametrelerin tüm olası değerleri için her bir zaman adımında hesaplamaları gerektirdiğinden, hesaplama karmaşıklığı olarak üsseldir. Hesaplama karmaşıklığını azaltmak için alfa-beta algoritmasının geri yönlü mesajlarının tümünü hesaplamak yerine, T anındaki sonsal dağılımdan, T zamanı için değişim noktası örneklenir (r_T). Örnek değeri için beta mesajını ($\beta_{t-1|t-1}$) hesaplayıp $t-1$ anı için sonsal dağılımları çıkarılır. Örnekleme bu şekilde sonsal dağılımlardan devam edilir. Bu ufak değişiklikle hesaplama karmaşıklığı geri yönde doğrusal olur.

Algoritmanın detayları aşağıdaki sözde kodla açıklanabilir:



Şekil 3: Zaman dilimlerine göre bağlantı parametrelerinin değişimi

Algoritma 1: İLERİ FİLTRELEME-GERİ ÖRNEKLEME(.)

```

for each  $p \in \{1, \dots, N^u\}$ 
  do  $C_p^u \sim Multn(\pi^u)$  örnekle
for each  $q \in \{1, \dots, N^s\}$ 
  do  $C_q^s \sim Multn(\pi^s)$  örnekle
for each  $t \in \{1, \dots, T\}$ 
  do  $r_t \sim Be(p)$  örnekle

```

Gibbs Örnekleme:

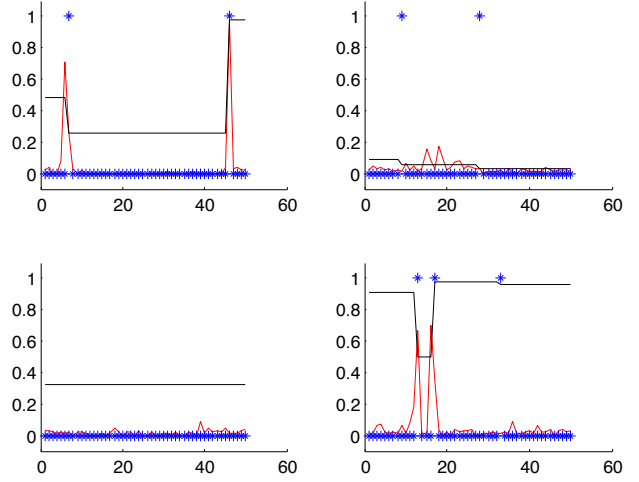
```

for each  $i \in \{1, \dots, epoch\}$ 
  for each  $t \in \{1, \dots, T\}$ 
    for each  $k \in \{1, \dots, K^u\}$ 
      do for each  $l \in \{1, \dots, K^s\}$ 
        do  $A_{t,kl}^i = k$  ve  $l$  kategorilerine dahil
           düğümlerden oluşan alt çizge
           İleri Süzgeçleme :
           for each  $A_{t,kl}^i$ 
             do  $\alpha_{t,kl}^i(r_{t,kl})$  mesajlarını hesapla
           Geri Örnekleme:
           for each  $t \in \{T, T-1, \dots, 1\}$ 
             do  $\begin{cases} s_t : \text{düzeltilmiş sonsal olasılık.} \\ r_t \sim s_t \text{ t dilimi değişim nokt.} \\ \beta_t(r_t) \text{ beta potansiyelleri} \end{cases}$ 
    for each  $p \in \{1, \dots, N^u\}$ 
      do  $C_p^u \sim p(C_p^u | Y_{1:T}, C_{p-}^u, C^s)$  örnekle
    for each  $q \in \{1, \dots, N^s\}$ 
      do  $C_q^s \sim p(C_q^s | Y_{1:T}, C_{q-}^s, C^u)$  örnekle

```

3. SONUÇLAR

Önerilen yeni algoritmayı test etmek için daha önce belirlenen zaman serisi modeline göre 50 zaman dilimi için veri üretildi. Üretici modelde 2 kullanıcı kategorisi ve 2 servis kategorisi olduğu varsayıldı. Her bir zaman diliminde



Şekil 4: Örneklenen değişim noktalarının sonsal dağılımları. Düz çizgiler bağlantı parametrelerinin değişimini, mavi yıldızlar gerçek değişim noktalarını, kırmızı grafikler de çıkarılan sonsal dağılımları gösteriyor.

değişim noktası olma olasılığı 0.04 olarak alındı. Her bir düğümün herhangi bir kategoriye ait olma olasılığı 0.5 olarak kabul edildi. Verideki gerçek değişim noktalarını, düğümlerin kategorilerini ve kategoriler arası bağlantı parametreleri şekil 3'te verilmiştir.

Önerilen algoritmanın ürettiği sonuçlar da şekil 4'de görülmektedir. Bu şekilden de anlaşıldığı gibi, Gibbs' örnekleme içinde geri yönlü örnekleme elde edilen noktalar, değişim noktaları için sonsal dağılım oluşturmaktadır. Gerçek değişim noktaları da, bu sonsal dağılımın en yüksek olasılık değerleriyle örtüşmektedir. Böylelikle önerilen algoritmanın daha az işlem karmaşıklığıyla tutarlı çıkarımlar yaptığı gözlemlendi.

4. KAYNAKÇA

- [1] Goldenberg, A. and Zheng, A. X. and Fienberg, S. E. and Airolidi, E. M., "A survey of statistical network models", arXiv, 2009, 6481222.
- [2] Airolidi, Edoardo M. and Blei, David M. and Fienberg, Stephen E. and Xing, Eric P., "Combining stochastic block models and mixed membership for statistical network analysis", ICML'06, 57-74, 2007.
- [3] Kurt, Barış and Cemgil, A. T., "Rastlantısız Öbek Çizgiler İçin Bayeşçi Model Seçimi", SIU 2011.
- [4] Bishop C., "Pattern Recognition and Machine Learning (Information Science and Statistics)", Springer, 2007.