

Ödev 1

Soru 1

24.05.2021 'de (başlık, 03.01.2023) github'da "fake news detection" adında veri madenciliği projesi paylaşılmış. Bu projede sahte haber tanıma sistemi yapılmıştır.

- Veri işlenmesi
- Veri temizlenmesi
- Veri dönüştürme
- Makine öğrenmesi algoritmaları

tekniklerini kullanarak spChack kullanıcısı adlı bir araştırmacı tarafından Jupyter Notebook'da yapılmıştır. (Logistic Regression, Naive Bayes, SVM ve Random Forests)

Soru 2

a) Veri madenciliği yaptığımız verilere genellikle bir veritabanından, veya hazırlanmış veri toplayıcı sensörler tarafından toplanır veya ulşunur. Bu veri toplanması veya veriye ulaşma aşamasındaki, sensör hataları, optiksel giriş hataları kullanıcısı hataları, sistem hataları yüzünden veriler ve bazı nitelikleri değişebilir.

Bir niteliğin durumu başka bir niteliğe bağlı olabilir.

Bir niteliğin uygunluğu diğer niteliğe bağlı olabilir. Örneğin şovanelerin sabit geliri olmaz. İnsanlar tablosu düştükçe ve sabit gelir adında bir nitelik belirtmek. Şovaneler için bu alan boş bırakılmak zorunda olacaktır.

b) Eksik nitelik değerlerinin önüne geçmek için, **elme yöntemi**, **eksik değeri tahmin etme**, analiz sırasında **görmezden gelme** ve olabilecek bir değeri **0** doldurma tekniklerini kullanabiliriz.

- Eksik değeri olan niteliğe sahip tüm verileri sileriz, elimine ederiz. Bu çok büyük veri setlerinde yapmak mantıklı olur çünkü analizdeki veri setini deho da küçültmek istemeyiz.

- Eksik değeri olan niteliğe sahip verilerin oranını, diğer verilere bakarak ortalamayı hesaplayarak bir veri girebiliriz. Bu durumda hem veri setimizi aynı boyuttaki bir veri de ortalamasını aldığımız için analiz aşamasında büyük, doğru bir seti yaratırız.

- Veriyi görmezden gelebiliriz kadar küçük katsayılar veririz.

Soru 3

Çok fazla öz nitelikli olan bir veri setinde öz nitelikleri azaltma aşamasında şunlar yapılabilir:

Missing Value Poth (Eksik değeri olan), belli bir seriye göre deho fazla eksik değeri içeren sütunların, analiz sırasında büyük farklar yaratmayacağı, bilinir. Bu nedenle bu sütunları sileriz.

Önceki metoda benzer olarak, hierarşide çok az değişim olan olan sütunlar analiz sırasında büyük değisikliklere yol almazlar. Bu nedenle bu sütunları da sileriz, görmezden geliriz.

Diğer bir farklı teknik ise Backward Feature Elimination (Geriye dönük silme) metodudur. Bu metod önce tüm değişkenleri siler. Ardından veri setinden sütunları silerek silme yapmaya devam eder. Enle edilen değisikliğe göre sütunun silinmesini veya katma değişkeni hesaplar. Bu şekilde katma değişkeni düşük olan sütunların silinmesi mantıklıdır.