

VERİ MADENCİLİĞİ

Veri Önleme

Yrd. Doç. Dr. Şule Gündüz Öğüdücü
<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

1

Konular

- Veri
- Veri Önleme
- Veriyi Tanıma
- Veri temizleme
- Veri birleştirme
- Veri dönüşümü
- Veri azaltma
- Benzerlik ve farklılık

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

2

Veri Nedir?

- nesneler ve nesnelerin niteliklerinden oluşan küme
 - kayıt (record), varlık (entity), örnek (sample, instance) nesne için kullanılabilir.
- Nitelik (attribute) bir nesnenin (object) bir özelliğidir
 - bir insanın yaşı, ortamın sıcaklığı..
 - boyut (dimension), özellik (feature, characteristic) olarak da kullanılır
- Nitelikler ve bu niteliklere ait değerler bir nesneyi oluşturur.

Tid	nitelikler				Dolan dırıcı
	Geri Odeme	Medeni Durum	Gelir		
1	Evet	Bekar	125K	-1	
2	Hayır	Evli	100K	-1	
3	Hayır	Bekar	70K	-1	
4	Evet	Evli	120K	-1	
5	Hayır	Boşanmış	95K	1	
6	Hayır	Evli	60K	-1	
7	Evet	Boşanmış	220K	-1	
8	Hayır	Bekar	85K	1	
9	Hayır	Evli	75K	-1	
10	Hayır	Bekar	90K	1	

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

3

Değer Kümeleri

- Nitelik için saptanmış sayılar veya semboller
- Nitelik & Değer Kümeleri
 - aynı nitelik farklı değer kümelerinden değer alabilir
 - ağırlık: kg, lb
 - farklı nitelikler aynı değer kümesinden değer alabilirler
 - ID, yaş: her ikisi de sayısal

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

4

Nitelik Türleri

- Belli aralıkta yeralan değişkenler (interval)
 - sıcaklık, tarih
- İkili değişkenler (binary)
 - cinsiyet
- Ayrık ve sıralı değişkenler (nominal, ordinal, ratio scaled)
 - göz rengi, posta kodu

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

5

Konular

- Veri
- Veri Önleme
- Veriyi Tanıma
- Veri temizleme
- Veri birleştirme
- Veri dönüşümü
- Veri azaltma
- Benzerlik ve farklılık

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

6

Problem

- Gerçek uygulamalarda toplanan veri kirlidir
 - eksik: bazı nitelik değerleri bazı nesneler için girilmemiş, veri madenciliği uygulaması için gerekli bir nitelik kaydedilmemiş
 - meslek = ""
 - gürültülü: hatalar var
 - maaş = "-10"
 - tutarsız: nitelik değerleri veya nitelik isimleri uyumsuz
 - yaş = "35", d.tarihi: "03/10/2004"
 - önceki oylama değerleri: "1,2,3", yeni oylama değerleri: "A,B,C"
 - bir kaynaktan nitelik değeri 'ad', diğerinde 'isim'

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

7

Verinin Gürültülü Olma Nedenleri

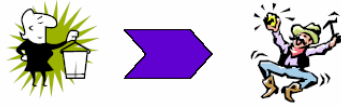
- Eksik veri kayıtlarının nedenleri
 - Veri toplandığı sırada bir nitelik değerinin elde edilememesi, bilinmemesi
 - Veri toplandığı sırada bazı niteliklerin gerekliliğinin görülememesi
 - İnsan, yazılım ya da donanım problemleri
- Gürültülü (hatalı) veri kayıtlarının nedenleri
 - Hatalı veri toplama gereçleri
 - İnsan, yazılım ya da donanım problemleri
 - Veri iletimi sırasında problemler
- Tutarsız veri kayıtlarının nedenleri
 - Verinin farklı veri kaynaklarında tutulması
 - İşlevsel bağımlılık kurallarına uyulmaması

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

8

Sonuç

- Veri güvenilirmez
 - Veri madenciliği sonuçlarına güvenilebilir mi?
 - Kullanılabilir veri madenciliği sonuçları kaliteli veri ile elde edilebilir
- Veri kaliteli ise veri madenciliği uygulamaları ile yararlı bilgi bulma şansı daha fazla.



<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

9

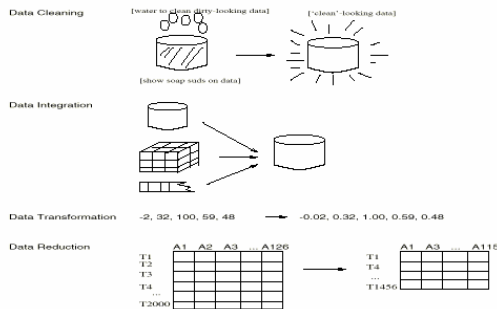
Veri Önışleme

- Veri temizleme
 - Eksik nitelik değerlerini tamamlama, hatalı veriyi düzeltme, aykırılıkları saptama ve temizleme, tutarsızlıkları giderme
- Veri birleştirme
 - Farklı veri kaynağındaki verileri birleştirme
- Veri dönüşümü
 - Normalizasyon ve biriktirme
- Veri azaltma
 - Aynı veri madenciliği sonuçları elde edilecek şekilde veri miktarını azaltma

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

10

Veri Önışleme



<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

11

Konular

- Veri
- Veri Önışleme
- Veriyi Tanıma
- Veri temizleme
- Veri birleştirme
- Veri dönüşümü
- Veri azaltma
- Benzerlik ve farklılık

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

12

Veriyi Tanımlayıcı Özellikler

- Amaç: Veriyi daha iyi anlamak
 - Merkezi eğilim (central tendency), varyasyon, yayılma, dağılım
- Verinin dağılım özellikleri
 - Ortanca, en büyük, en küçük, sıklık derecesi, aykırılık, varyans
- Sayısal nitelikler -> sıralanabilir değerler
 - verinin dağılımı
 - kutu grafiği çizimi ve sıklık derecesi incelemesi

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

13

Merkezi Eğilimi Ölçme

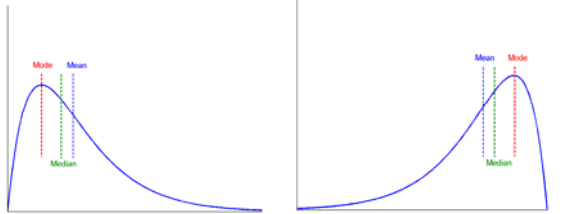
- Ortalama:
 - ağırlıklı ortalama
 - kırılmış ortalama: Uç değerleri kullanmadan hesaplama
- Ortanca (median): Verinin tümü kullanılarak hesaplanır
 - veri sayısı tek ise ortadaki değer, çift sayı ise ortadaki iki değer ortalaması
- Mod
 - Veri içinde en sıklıkla görülen değer
 - Unimodal, bimodal, trimodal
 - deneysel formül: $mean - mode = 3 \times (mean - median)$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

14

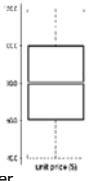
Simetrik – Asimetrik Dağılım

- Simetrik ve asimetrik verinin ortalama, ortanca ve mod değerleri



Verinin Dağılımını Ölçme

- Çeyrek, aykırılıklar, kutu grafiği çizimi
 - Çeyrek (quartile) : nitelik değerleri küçükten büyüğe doğru sıralanır.
 - Q1: ilk %25, Q3: ilk %75
 - Dörtü aralık (Inter-quartile Range): IQR= Q3-Q1
 - Five Number Summary: min, Q1, median, Q3, max
 - Kutu Grafiği Çizimi:
 - Q1 ve Q3 aralığında bir kutu
 - kutu içinde ortanca noktayı gösteren bir çizgi
 - kutudan min ve max değerlere birer uzantı
 - Aykırılıklar: 1,5xIQR değerinden küçük/büyük olan değerler
- Varyans ve standart sapma



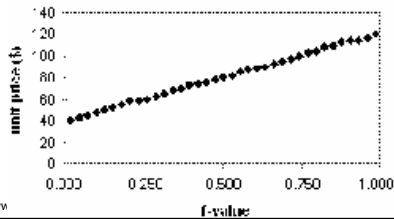
$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2 \right] \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^n (x_i - \mu)^2 = \frac{1}{N} \sum_{i=1}^n x_i^2 - \mu^2$$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

16

Quantile Plot

- Verinin bütünü bir nitelik değerine göre görüntüleme
- Veri bir nitelik değerine göre küçükten büyüğe doğru sıralanır
 - x_i değeri için % 100 f_i miktardaki veri x_i değerinden küçük ya da eşittir

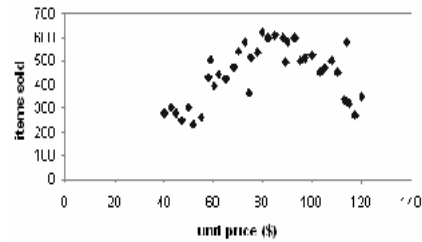


<http://www>

17

Scatter Plot

- İki sayısal nitelik değeri arasındaki ilişkiyi görmek



<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

18

Konular

- Veri
- Veri Önışleme
- Veriyi Tanıma
- Veri temizleme
- Veri birleřtirme
- Veri dönüşümü
- Veri azaltma
- Benzerlik ve farklılık

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

19

Veri Temizleme

- Gerçek uygulamalarda veri eksik, gürültülü veya tutarsız olabilir.
- Veri temizleme işlemleri
 - Eksik nitelik değerlerini tamamlama
 - Aykırılıkların bulunması ve gürültülü verinin düzeltilmesi
 - Tutarsızlıkların giderilmesi



<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

20

Eksik Veri

- Veri için bazı niteliklerin değerleri her zaman bilinemeyebilir.
- Eksik veri
 - diğer veri kayıtlarıyla tutarsızlığı nedeniyle silinmesi
 - Bazı nitelik değerleri hatalı olması dolayısıyla silinmesi
 - yanlış anlama sonucu kaydedilmeme
 - veri giriři sırasında bazı nitelikleri önemsiz görme

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

21

Eksik Veriler nasıl Tamamlanır?

- Eksik nitelik değerleri olan veri kayıtlarını kullanma
- Eksik nitelik değerlerini elle doldur
- Eksik nitelik değerleri için global bir değışken kullan (Null, bilinmiyor,...)
- Eksik nitelik değerlerini o niteliğin ortalama değeri ile doldur
- Aynı sınıfa ait kayıtların nitelik değerlerinin ortalaması ile doldur
- Olasılığı en fazla olan nitelik değeriyle doldur

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

22

Gürültülü Veri

- Ölçülen bir değerdeki hata
- Yanlış nitelik değerleri
 - hatalı veri toplama gereçleri
 - veri giriři problemleri
 - veri iletimi problemleri
 - teknolojik kısıtlar
 - nitelik isimlerinde tutarsızlık

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

23

Gürültülü Veri nasıl Düzeltir?

- Gürültüyü yok etme
 - Bölmeleme
 - veri sıralanır, eşit genişlik veya eşit derinlik ile bölünür
 - Demetleme
 - aykırılıkları belirler
 - Eğri uydurma
 - veriyi bir fonksiyona uydurarak gürültüyü düzeltir
- <http://control.cs.berkeley.edu/abc/>

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

24

Bölmeleme

- Veri sıralanır: 4, 8, 15, 21, 21, 24, 25, 28, 34
 - Eşit genişlik: Bölme sayısı belirlenir. Eşit aralıklarla bölünür
 - Eşit derinlik: Her bölmede eşit sayıda örnek kalacak şekilde bölünür.
 - her bölme ortalamayla ya da bölmenin en alt ve üst sınırlarıyla temsil edilir

Bölme genişliği:3
 1. Bölme: 4, 8, 15
 2. Bölme: 21, 21, 24
 3. Bölme: 25, 28, 34

Ortalamayla düzeltme:
 1. Bölme: 9, 9, 9
 2. Bölme: 22, 22, 22
 3. Bölme: 29, 29, 29

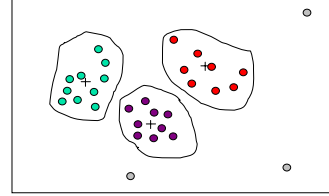
Alt-üst sınırla düzeltme:
 1. Bölme: 4, 4, 15
 2. Bölme: 21, 21, 24
 3. Bölme: 25, 25, 34

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

25

Demetleme

- Benzer veriler aynı demette olacak şekilde gruplanır
- Bu demetlerin dışında kalan veriler aykırılık olarak belirlenir ve silinir

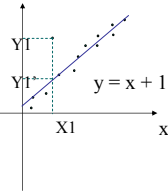


<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

26

Eğri Uydurma

- Veri bir fonksiyona uydurulur. Doğrusal eğri uydurmada, bir değişkenin değeri diğer bir değişken kullanılarak bulunabilir.



<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

27

Konular

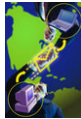
- Veri
- Veri Önışleme
- Veriyi Tanıma
- Veri temizleme
- Veri birleştirme
- Veri dönüşümü
- Veri azaltma
- Benzerlik ve farklılık

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

28

Veri Birleştirme

- Farklı kaynaklardan verilerin tutarlı olarak birleştirilmesi
- Şema birleştirilmesi
 - Aynı varlıkların saptanması: A.cust_id=B.cust_num
 - meta veri kullanılır
- Nitelik değerlerinin tutarsızlığının saptanması
 - Aynı nitelik için farklı kaynaklarda farklı değerler olması
 - Farklı metrikler kullanılması



<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

29

Gereksiz Veri

- Farklı veri kaynaklarından veriler birleştirilince gereksiz (fazla) veri oluşabilir

- aynı nitelik farklı kaynaklarda farklı isimle
- bir niteliğin değeri başka bir nitelik kullanılarak hesaplanabilir
 - korelasyon hesaplaması: sayısal nitelikler
 - =0: nitelikler bağımsız, >0: pozitif korelasyon, <0: negatif korelasyon

$$r_{A,B} = \frac{\sum (A - \bar{A})(B - \bar{B})}{(n-1)\sigma_A\sigma_B} \quad \bar{A} = \frac{\sum A}{n} \quad \sigma_A = \sqrt{\frac{\sum (A - \bar{A})^2}{n-1}}$$

- korelasyon hesaplaması: ayrık nitelikler (chi-square test)

$$\chi^2 = \sum \frac{(Observed - Expected)^2}{Expected}$$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

30

Konular

- Veri
- Veri Önışleme
- Veriyi Tanıma
- Veri temizleme
- Veri birleřtirme
- Veri dönüşümü
- Veri azaltma
- Benzerlik ve farklılık

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

31

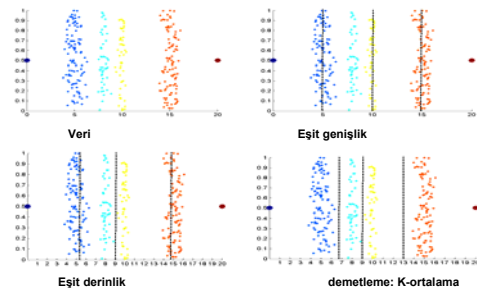
Veri Dönüşümü

- Veri, veri madencilięi uygulamaları için uygun olmayabilir
 - Seçilen algoritmaya uygun olmayabilir
 - Veri belirleyici deęil
- Çözüm
 - Veri düzeltme
 - Bölmeleme
 - Demetleme
 - Eğri Uydurma
 - Biriktirme
 - Genelleme
 - Normalizasyon
 - Nitelik oluřturma

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

32

Bölmeleme ve Demetleme Yöntemleri ile Veri Düzeltme



<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

33

Normalizasyon

- min-max normalizasyon

$$v' = \frac{v - \min_A}{\max_A - \min_A} (\text{new_max}_A - \text{new_min}_A) + \text{new_min}_A$$

- z-score normalizasyon

$$v' = \frac{v - \text{mean}_A}{\text{stand_dev}_A}$$

- ondalık normalizasyon

$$v' = \frac{v}{10^j} \quad j: \text{Max}(|v'|) < 1 \text{ olacak řekildeki en küçük tam sayı}$$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

34

Nitelik Oluřturma

- Yeni nitelikler yarat
 - orjinal niteliklerden daha önemli bilgi içersin
 - alan=boy x en
 - veri madencilięi algoritmalarının başarımı daha iyi olsun

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

35

Konular

- Veri
- Veri Önışleme
- Veriyi Tanıma
- Veri temizleme
- Veri birleřtirme
- Veri dönüşümü
- Veri azaltma
- Benzerlik ve farklılık

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

36

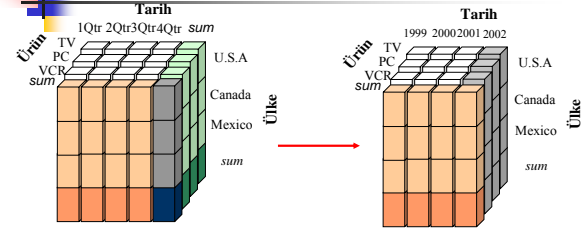
Veri Azaltma

- Veri miktarı çok fazla olduğu zaman veri madenciliği algoritmalarının çalışması ve sonuç üretmesi çok uzun sürebilir
 - veriyi azaltma başarımı artırır
 - sonucun (nerdeyse) hiç değişmemesi gerekir
- Veri azaltma
 - nitelik birleştirme
 - nitelik azaltma
 - veri sıkıştırma
 - veri ayrıştırma ve kavram oluşturma
 - veri küçültme
 - eğri uydurma
 - demetleme
 - histogram
 - örnekleme

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

37

Nitelik Birleştirme



- Sorgulama için gerekli olan boyutlar kullanılıyor.

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

38

Nitelik Seçme - Nitelik Azaltma

- Nitelik Seçme
 - Nitelikler kümesinin bir alt kümesi seçilerek veri madenciliği işlemi yapılır
- Nitelik azaltma
 - d boyutlu veri kümesi $k < d$ olacak şekilde k boyuta taşınır

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

39

Nitelik Seçme

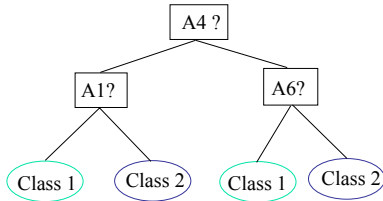
- Nitelik seçme
 - Veri madenciliği uygulaması için gerekli olan niteliklerin seçilmesi
 - Nitelikler altkümesi kullanılarak elde edilen sınıfların dağılımları gerçek dağılıma eşit ya da çok yakın olmalı
 - Veri madenciliği işlemi yer ve zaman karmaşıklığını azaltma
 - Sistemin başarımını artırma
- Sezgisel yöntemler kullanılarak nitelikler seçilebilir.
 - istatistiksel anlamlılık testi (statistical significance)
 - bilgi kazancı (information gain)
 - karar ağaçları

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

40

Örnek

Başlangıç nitelikler kümesi:
 $\{A1, A2, A3, A4, A5, A6\}$



Seçilen nitelik kümesi: $\{A1, A4, A6\}$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

41

Nitelik Azaltma

- Çok boyutlu veriyi daha küçük boyutlu uzaya taşıma
- d nitelikten oluşan n adet veri $D = \{x_1, x_2, \dots, x_n\}$ k boyutlu uzaya taşınır:

$$x_i \in \mathbb{R}^d \rightarrow y_i \in \mathbb{R}^k (k \ll d)$$

- Veri kümesinde yer alan bütün nitelikler kullanılır
 - Niteliklerin doğrusal kombinasyonu
- Niteliklerin ayırtıcılığına artırma

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

42

Veri Sıkıştırma

- Verinin boyutunu azaltır
 - daha az saklama ortamı
 - veriye ulaşmak daha çabuk
- Kayıplı ve kayıpsız veri sıkıştırma
 - bazı yöntemler bazı veri tiplerine uygun
 - her veri tipi için kullanılan yöntemler de var
- Eğer veri madenciliği yöntemi sıkıştırılmış veri üzerinde doğrudan çalışabiliyorsa elverişli

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

43

Veri Ayrıştırma

- Bazı veri madenciliği algoritmaları sadece ayrıık veriler ile çalışır
- Sürekli bir nitelik değerini bölerek her aralığı etiketler
- Verinin değeri, bulunduğu aralığın etiketi ile değişir
- Veri boyutu küçülür
- Kavram oluşturmak için kullanılır

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

44

Kavram Oluşturma

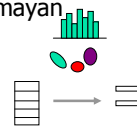
- Sayısal veriler
 - çok geniş aralıkta olabilir
 - değerleri çok sık değişebilir
- Sayısal veriler için kavram oluşturma
 - bölmeleme
 - histogram
 - demetleme
 - entropi

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

45

Veri Küçültme

- Veriyi farklı şekillerde gösterme
 - parametrik
 - eğri uydurma
 - parametrik olmayan
 - histogram
 - demetleme
 - örnekleme



<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

46

Histogram ile Veri Küçültme

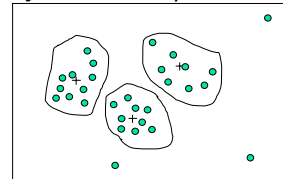
- Verinin dağılımı
- Veriyi bölerek her bölüm için veri değerini gösterir (toplam, ortalama)
 - eşit genişlik (equi-width): bölmelerin genişliği eşit
 - eşit yükseklik (equi-height): her bölmedeki veri sayısı eşit
 - v-optimal: en az varyansı olan histogram $\Sigma(\text{count}_i * \text{value}_i)$
 - MaxDiff: bölme genişliğini kullanıcı belirler

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

47

Demetleme ile Veri Küçültme

- Veri demetlere ayrılır
- Veri demetleri temsil eden örnekler (demet merkezleri) ve aykırılıklar ile temsil edilir
- Etkisi verinin dağılımına bağlı
- Hiyerarşik demetleme yöntemleri kullanılabilir.



<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

48

Örnekleme ile Veri Küçültme

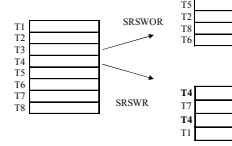
- Büyük veri kümesini daha küçük bir alt küme ile temsil etme
- Alt küme nasıl seçiliyor?
 - yerine koymadan örnekleme (SRSWOR)
 - yerine koyarak örnekleme (SRSWR)
 - demet örnekleme (yerine koymadan veya koyarak)
 - katman örnekleme (katman: nitelik değerine göre grup)

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

49

Örnek

- Örnekleme

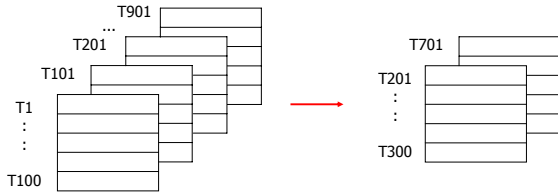


<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

50

Örnek

- Demetleme



<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

51

Örnek

- Katman Örnekleme

T1	genç	T1	genç
T2	genç	T4	genç
T3	genç	T6	orta yaşlı
T4	genç	T7	orta yaşlı
T5	orta yaşlı	T9	orta yaşlı
T6	orta yaşlı	T11	orta yaşlı
T7	orta yaşlı	T13	yaşlı
T8	orta yaşlı		
T9	orta yaşlı		
T10	orta yaşlı		
T11	orta yaşlı		
T12	orta yaşlı		
T13	yaşlı		
T14	yaşlı		

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

52

Konular

- Veri
- Veri Önışleme
- Veriyi Tanıma
- Veri temizleme
- Veri birleřtirme
- Veri dönüşümü
- Veri azaltma
- Benzerlik ve farklılık

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

53

Benzerlik ve Farklılık

- Benzerlik
 - iki nesnenin benzerliğini ölçen sayısal değer
 - nesneler birbirine daha benzer ise daha büyük
 - genelde 0-1 aralığında değer alır
- Farklılık
 - iki nesnenin birbirinden ne kadar farklı olduğunu gösteren sayısal değer
 - nesneler birbirine daha benzer ise daha küçük
 - en küçük farklılık genelde 0
 - üst sınır değişebilir

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

54

Öklid Uzaklığı

Veri kümesi

$$\begin{bmatrix} x_{11} & \dots & x_{1f} & \dots & x_{1p} \\ \dots & \dots & \dots & \dots & \dots \\ x_{i1} & \dots & x_{if} & \dots & x_{ip} \\ \dots & \dots & \dots & \dots & \dots \\ x_{n1} & \dots & x_{nf} & \dots & x_{np} \end{bmatrix}$$

Uzaklık matrisi

$$\begin{bmatrix} 0 & & & & \\ d(2,1) & 0 & & & \\ d(3,1) & d(3,2) & 0 & & \\ \vdots & \vdots & \vdots & \ddots & \\ d(n,1) & d(n,2) & \dots & \dots & 0 \end{bmatrix}$$

- Öklid uzaklığı (Euclidean Distance) nesneler arası farklılığı bulmak için kullanılır.

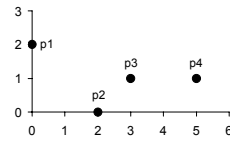
- p adet niteliği (boyutu) olan i ve j nesneleri arasındaki uzaklık

$$d(i, j) = \sqrt{|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2}$$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

55

Örnek: Öklid Uzaklığı



nesne	x	y
p1	0	2
p2	2	0
p3	3	1
p4	5	1

	p1	p2	p3	p4
p1	0	2.828	3.162	5.099
p2	2.828	0	1.414	3.162
p3	3.162	1.414	0	2
p4	5.099	3.162	2	0

Uzaklık Matrisi

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

56

Minkowski Uzaklığı

- Öklid uzaklığının genelleştirilmiş hali

$$d(i, j) = \sqrt[q]{(|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots + |x_{ip} - x_{jp}|^q)} \quad q: \text{pozitif tam sayı}$$

- $q=1 \rightarrow$ Manhattan uzaklığı

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

57

Uzaklık Özellikleri

- $q=1 \Rightarrow$ Manhattan Uzaklığı
- $q=2 \Rightarrow$ Öklid Uzaklığı
- Uzaklık ölçütünün sağlaması gereken özellikler:
 - $d(i, j) \geq 0$
 - $d(i, i) = 0$
 - $d(i, j) = d(j, i)$
 - $d(i, j) \leq d(i, h) + d(h, j)$
- Uzaklıklar ağırlıklı olarak da hesaplanabilir:

$$d(i, j) = \sqrt{w_1 |x_{i1} - x_{j1}|^2 + w_2 |x_{i2} - x_{j2}|^2 + \dots + w_p |x_{ip} - x_{jp}|^2}$$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

58

Benzerlik Özellikleri

- İki nesne arası benzerlik özellikleri

- $\text{sim}(i, j) \geq 0$
- $\text{sim}(i, j) = \text{sim}(j, i)$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

59

İkili Değişkenler Arası Benzerlik

- İkili bir değişkenin 0 veya 1 olarak iki değeri olabilir.
- Bir olasılık tablosu oluşturulur:

Nesne i	Nesne j	
	0	1
0	M_{00}	M_{01}
1	M_{10}	M_{11}

M_{00} : i nesnesinin 0, j nesnesinin 0 olduğu niteliklerin sayısı
 M_{01} : i nesnesinin 1, j nesnesinin 0 olduğu niteliklerin sayısı
 M_{10} : i nesnesinin 0, j nesnesinin 1 olduğu niteliklerin sayısı
 M_{11} : i nesnesinin 1, j nesnesinin 1 olduğu niteliklerin sayısı

- Yalın uyum katsayısı (simple matching coefficient): ikili değişkenin simetrik olduğu durumlarda

$$\text{sim}(i, j) = \frac{M_{11} + M_{00}}{M_{00} + M_{01} + M_{10} + M_{11}}$$

- Jaccard katsayısı (İkili değişkenin asimetrik olduğu durumlar):

$$d(i, j) = \frac{M_{11}}{M_{01} + M_{10} + M_{11}}$$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?eId=195>

60

Örnek

$p = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0\ 0$
 $q = 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 0\ 1$

$M_{01} = 2$
 $M_{10} = 1$
 $M_{00} = 7$
 $M_{11} = 0$

Yalın Uyum Katsayısı:

$$(M_{11} + M_{00}) / (M_{01} + M_{10} + M_{11} + M_{00}) = (0+7) / (2+1+0+7) = 0.7$$

Jaccard Katsayısı:

$$(M_{11}) / (M_{01} + M_{10} + M_{11}) = 0 / (2 + 1 + 0) = 0$$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?elId=195>

61

Kosinüs Benzerliği

- d_1 ve d_2 iki doküman. Kosinüs benzerliği
 $\cos(d_1, d_2) = d_1 \bullet d_2 / ||d_1|| \cdot ||d_2||$
 $d_1 \bullet d_2$: iki dokümanın vektör çarpımı
 $||d_i||$: d_i dokümanının uzunluğu

- Örnek

$d_1 = 3\ 2\ 0\ 5\ 0\ 0\ 0\ 2\ 0\ 0$
 $d_2 = 1\ 0\ 0\ 0\ 0\ 0\ 0\ 1\ 0\ 2$

$$\begin{aligned} d_1 \bullet d_2 &= 3*1 + 2*0 + 0*0 + 5*0 + 0*0 + 0*0 + 0*0 + 2*1 + 0*0 + 0*2 = 5 \\ ||d_1|| &= (3^2 + 2^2 + 0^2 + 5^2 + 0^2 + 0^2 + 0^2 + 2^2 + 0^2 + 0^2)^{0.5} = (42)^{0.5} = 6.481 \\ ||d_2|| &= (1^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 0^2 + 1^2 + 0^2 + 2^2)^{0.5} = (6)^{0.5} = 2.245 \\ \cos(d_1, d_2) &= .3150 \end{aligned}$$

<http://www.ninova.itu.edu.tr/EgitimDetay.aspx?elId=195>

62