

METİN MADENCİLİĞİ

Hazırlayan:

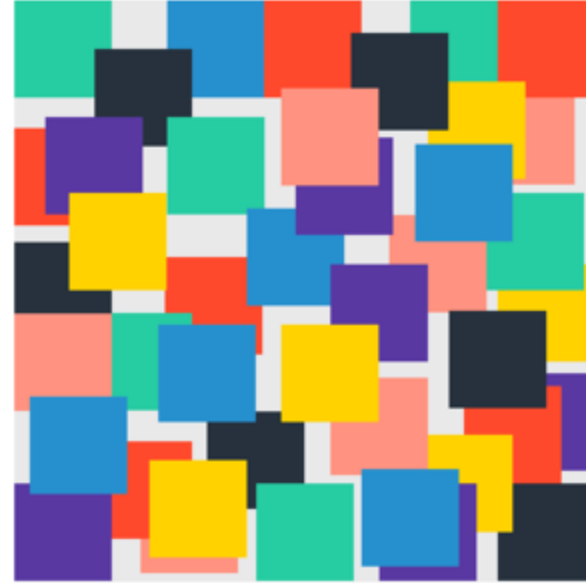
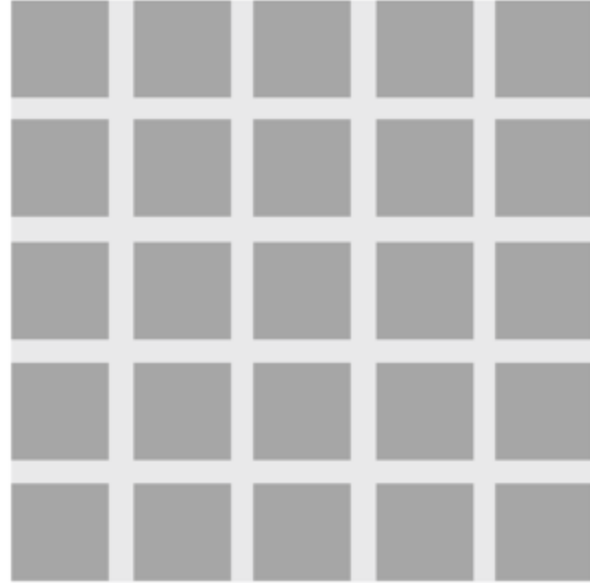
Dr. Ekin EKİNCİ



METİN MADENCİLİĞİNE GENEL BAKIŞ (1)*

Yapısal Veri

Yapısal Olmayan Veri



Veri tabanı, CRM, ERP

Metin, Ses, Video

Bugün üretilen tüm verilerin yüzde 80'inden fazlası yapısal olmayan veri olarak kabul edilmektedir.

METİN MADENCİLİĞİNE GENEL BAKIŞ (2)

VERİ

Yapısal Veri

Yarı Yapısal Veri

Yapısal Olmayan Veri

İlan No	650143961
İlan Tarihi	11 Mayıs 2019
Marka	BMW
Seri	X5
Model	40e xDrive M Plus
Yıl	2017
Yakıt	Hybrid
Vites	Otomatik
KM	22.000
Kasa Tipi	SUV
Motor Gücü	313 hp
Motor Hacmi	1997 cc
Çekiş	4x4
Kapı	5
Renk	Gümüş Gri
Garanti	Evet
Plaka / Uyruk	Türkiye (TR) Plakalı
Kimden	Galeriden
Takas	Hayır
Durumu	İkinci El

```
{
  "person": {
    "name": "Jennifer",
    "surname": "Green",
    "friends": [
      {
        "age": 25,
        "isDeveloped": true
      },
      {
        "age": 28,
        "isDeveloped": false
      }
    ]
  },
  "salary": 1000
}
```

Hürriyet Lezizz @HurriyetLezizz · 8 sa.

Kavurma yapılışı en kolay ancak en lezzetli geleneksel lezzetlerimizden biri. İftarda yanında pilavla zevkle yiyebileceğiniz kuzu kavurma tarifimize bir göz atmaz mıydınız? hry.yt/sHi3Z

METİN MADENCİLİĞİNE GENEL BAKIŞ (3)*

	Structured Data	Unstructured Data
Characteristics	<ul style="list-style-type: none">• Pre-defined data models• Usually text only• Easy to search	<ul style="list-style-type: none">• No pre-defined data model• May be text, images, sound, video or other formats• Difficult to search
Resides in	<ul style="list-style-type: none">• Relational databases• Data warehouses	<ul style="list-style-type: none">• Applications• NoSQL databases• Data warehouses• Data lakes
Generated by	Humans or machines	Humans or machines
Typical applications	<ul style="list-style-type: none">• Airline reservation systems• Inventory control• CRM systems• ERP systems	<ul style="list-style-type: none">• Word processing• Presentation software• Email clients• Tools for viewing or editing media
Examples	<ul style="list-style-type: none">• Dates• Phone numbers• Social security numbers• Credit card numbers• Customer names• Addresses• Product names and numbers• Transaction information	<ul style="list-style-type: none">• Text files• Reports• Email messages• Audio files• Video files• Images• Surveillance imagery

METİN MADENCİLİĞİNE GENEL BAKIŞ (4)

- Yapısal ve yapısal olmayan milyarlarca içeriği biz kullanıcılarına sunan Web, günümüzün önemli veri kaynaklarından birisi haline gelmiştir.
- Sunulan içerik her geçen gün büyümektedir.
- İçeriğin %80'i dokümanlar şeklinde organize edilmiştir: haberler, forumlar, e-mailler, haber grupları, sosyal medya, ...

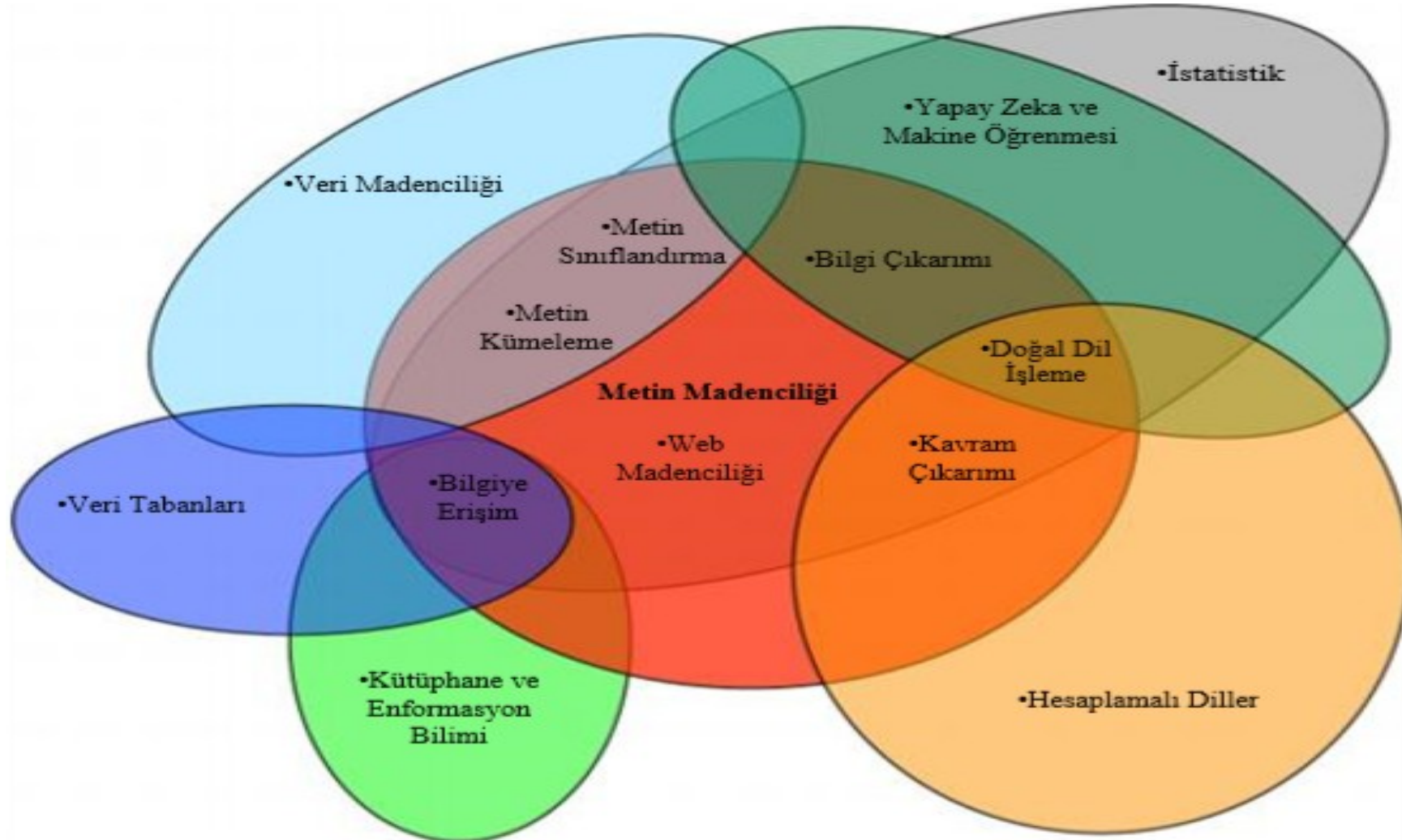
METİN MADENCİLİĞİNE GENEL BAKIŞ (5)

- Sunulan içerik her geçen gün büyümekte, bu içerikten **istenilen bilginin otomatik bir şekilde çıkartılması ve çıkartılan bilginin:**
 - **Organize edilme,**
 - **Analiz edilme ve**
 - **Anlaşılması adımıında ise metin madenciliğine ihtiyaç duyulmaktadır.**

METİN MADENCİLİĞİ NEDİR? (1)

- Her geçen gün artan veri miktarı bu verileri yönetmeyi ve içerisinden önemli olan ancak keşfedilmemiş bilgiyi çıkarmayı gerekli hale getirmiş ve metin madenciliği kavramı ortaya çıkmıştır.
- Eldeki dokümanlardan belli bir amaç çerçevesinde önceden bilinmeyen ancak potansiyel olarak faydalı bilginin çıkarılması şeklinde tanımlanmaktadır (Visa, 2001).

METİN MADENCİLİĞİ NEDİR? (2)*



METİN MADENCİLİĞİ BİLEŞENLERİ (1)

- **Bilgi Edinme (Information Retrieval):** Büyük koleksiyonlardan (genellikle bilgisayarlarda saklanan) belli bir amaca yönelik bir bilgi ihtiyacını karşılayan, yapılandırılmamış nitelikte bir materyalin (genellikle belgeler) elde edilmesidir (*).
- **Doküman Kümeleme (Document Clustering):** Büyük miktardaki doküman koleksiyonunu her birinin bir konuyu temsil ettiği az sayıdaki anlamlı kümelere dağıtma görevidir.
- **Doküman Sınıflandırma (Document Classification):** Dokümanlarını önceden tanımlı bir ya da daha fazla sınıfa atama görevidir.

METİN MADENCİLİĞİ BİLEŞENLERİ (2)

- **Web Madenciliği (Web Mining):** Veri madenciliğinin alt dallarından biri olan web madenciliği webden elde edilen verilerden bilginin çıkartılmasını amaçlar.
- **Bilgi Çıkarımı (Information Extraction):** Yapılandırılmamış dokümanlardan yapılandırılmış bilginin çıkartılması görevidir.
- **Doğal Dil İşleme (Natural Language Processing):** doğal dil üzerine inceleme, çözümleme, yorumlama, bilgi çıkarma, üretme yapan bilgisayar sistemi şeklinde tanımlanmaktadır (Oğuzlar, 2011).

METİN MADENCİLİĞİ BİLEŞENLERİ (3)

- **Kavram Çıkarımı (Concept Extraction):** Kelimelerin ve öbeklerin anlamsal olarak benzer gruplar altında öbeklenmesi görevidir.

METİN MADENCİLİĞİ UYGULAMA ALANLARI

- Konu çıkarımı
- Duygu analizi
- Soru cevaplama sistemleri
- Yazar analizi
- Doküman özetleme
- Haberlerin sınıflandırılması
- Spam filtreleme,...

METİN MADENCİLİĞİ ADIMLARI



METNİN ELDE EDİLMESİ

- Metin madenciliği adımlarını gerçekleştirebilmemiz için ilk olarak amaca yönelik bir veri kümesinin elde edilmesi gerekmektedir.
 - Hazır veri kümelerini kullanabiliriz:
 - UCI Machine Learning Repository
 - Kaggle
 - Kemik Doğal Dil İşleme Grubu
 - Kendi veri kümemizi kendimiz oluşturabiliriz:
 - Web Crawler ile

METİN ÖNİŞLEME

- Metinler üzerinde yapılacak önışleme alıřılacak amaca gre farklılıklar gstermekle birlikte temel önışleme adımları;
 - noktalama işareleri, sayı ve özel karakterlerin eldeki metinlerden ıkartılması,
 - büyük küçük harf duyarlı olmamasından tür büyük harflerin küçük harflere dnüştürölmesi,
 - metni meydana getiren ve ok sık tekrarlanan ancak doküman için önemli olmayan durak kelimelerinin eldeki metinlerden ayıklanması,
 - yazım hatalarının düzeltilmesi (normalizasyon),
 - POS tagging
 - gövdelemenin gereklenmesi řeklinde sıralanmaktadır.
- Bu önışleme adımları doęal dil işleme sürecini oluřurmaktadır.

METİN ÖNİŞLEME

Yorumun ilk hali

...There was far too much bread, not enough gelato or pudding... However, people like trendy spots because of celebrity sightings... today's celebrity sighting was Vivica_Fox, enjoying patio seating and a the company iof a couple of good friends.

LanguageTool kütüphanesi ile yazım hatalarının düzeltilmesi

(<http://wiki.languagetool.org/java-api>)



Yazım hatalarının düzeltilmesi

...There was far too much bread, not enough gelato or pudding... However, people like trendy spots because of celebrity sightings... today's celebrity sighting was Vivica_Fox, enjoying patio seating and a the company **of** a couple of good friends.

Stanford Üniversitesi'nin doğal dil işleme aracı

<http://nlp.stanford.edu/software/>



Morfolojik analiz

...bread enough relate pudding however people like trendy **spot** celebrity **sighting** today celebrity sighting vivica_fox **enjoy** patio **seat** company couple good **friend**

Yorumların düzenlenmiş hali

...bread enough relate pudding **however** people like trendy spots celebrity sightings **today** celebrity sighting **vivica_fox** enjoying patio seating company couple good friends

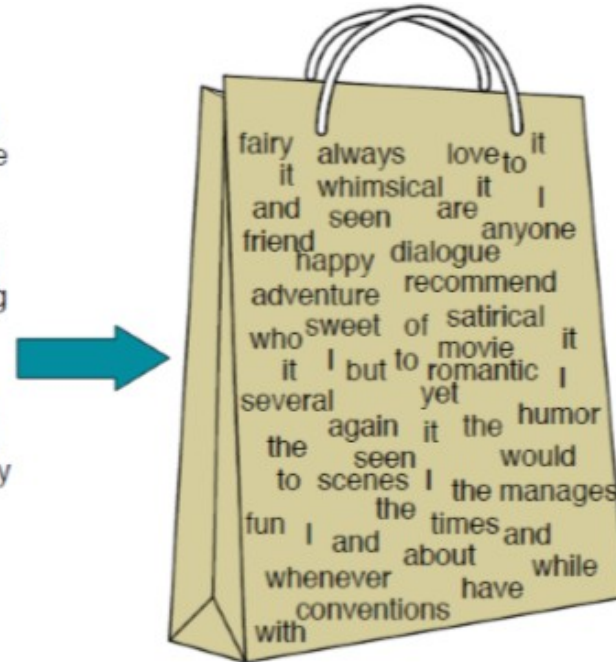
Noktalama işaretlerinin temizlenmesi, büyük küçük harf dönüşümü, durak kelimelerinin ayıklanması



METİN DÖNÜŞÜMÜ (1)

- **Kelime torbası:** Bir dokümanın tipik temsilidir. Kelimeler frekansları ile temsil edilmektedir, kelimelerin doküman içerisindeki konumu göz ardı edilmektedir.
- Kelimelerin ağırlıklarının hesaplanması gerekmektedir.
- Ağırlık hesabı ise kelimenin ilgili sınıfta geçme sıklığı şeklinde hesaplanır.

I love this movie! It's sweet, but with satirical humor. The dialogue is great and the adventure scenes are fun... It manages to be whimsical and romantic while laughing at the conventions of the fairy tale genre. I would recommend it to just about anyone. I've seen it several times, and I'm always happy to see it again whenever I have a friend who hasn't seen it yet!



it	6
I	5
the	4
to	3
and	3
seen	2
yet	1
would	1
whimsical	1
times	1
sweet	1
satirical	1
adventure	1
genre	1
fairy	1
humor	1
have	1
great	1
...	...

METİN DÖNÜŞÜMÜ (2)

- **Vektör Uzayı Modeli:** Dokümanların ortak bir uzayda vektörler olarak gösterilmesi, vektör uzay modeli olarak ifade edilmektedir.
- Bu modelde dokümanlar ağırlık vektörü olarak temsil edilmektedir.
- Terim ağırlıkları Tf (term frequency) ya da Tf-Idf (term frequency-invert document frequency) şemalarına göre hesaplanmaktadır.

METİN DÖNÜŞÜMÜ (3)

- **Terim frekansı:** Bir terimin ilgili dokümanda kaç kere geçtiğini temsil etmektedir.
- i . terimin j . dokümandaki frekansı ile temsil edilmektedir.

	Doc 1	Doc 2	...	Doc n
Term(s) 1	12	2	...	1
Term(s) 2	0	1	...	0
...
Term(s) n	0	6	...	3

METİN DÖNÜŞÜMÜ (4)

- 8 **Ters metin frekansı:** Eğer bir terim çok fazla sayıda dokümanda bulunuyorsa muhtemelen bu terim gerçekleştirilecek görev için önemli değildir.
- N toplam doküman sayısı iken df_i i. terimin geçtiği toplam doküman sayısıdır.

$$idf_i = \log \frac{N}{df_i}$$

METİN DÖNÜŞÜMÜ (5)

⊗ i. terimin j. dokümandaki ağırlığı $w_{i,j}$ ile temsil edilmektedir.

$$w_{i,j} = tf_{i,j} * idf_i$$

$$\begin{pmatrix} & D_1 & D_2 & \dots & D_t \\ T_1 & w_{11} & w_{21} & \dots & w_{t1} \\ T_2 & w_{12} & w_{22} & \dots & w_{t2} \\ \vdots & \vdots & \vdots & & \vdots \\ \vdots & \vdots & \vdots & & \vdots \\ T_n & w_{1n} & w_{2n} & \dots & w_{tn} \end{pmatrix}$$

METİN DÖNÜŞÜMÜ (5)

- **Cosinüs Benzerliği:** Cosinüs benzerliği, iki vektör arasındaki açının cosinüsünü ölçer.
- İki dokümanın ağırlık vektörleri üzerinden benzerliklerini ölçmek için cosinüs benzerliğinden yararlanılmaktadır.

$$\cos(d_j, d_k) = \frac{\langle d_j, d_k \rangle}{\|d_j\| \|d_k\|} = \frac{\sum_{i=1}^V w_{i,j} \times w_{i,k}}{\sqrt{\sum_{i=1}^V w_{i,j}^2} \sqrt{\sum_{i=1}^V w_{i,k}^2}}$$

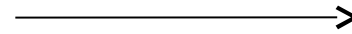
ÖZELLİK SEÇİMİ

- Model oluşturmada kullanılacak önemli özelliklerin bir alt kümesini seçme işlemidir.
- Gereksiz özellikler veri madenciliği görevi için herhangi bir katkı sağlamamaktadır.
- Ki-kare istatistiği, tekil değer ayrışımı, doküman frekansı için bir eşik değeri belirleme,...

VERİ MADENCİLİĞİ (1)

- **Sınıflandırma:** Dokümanın önceden tanımlanmış sınıflardan birine atanmasıdır.

w_{11}	w_{12}	w_{13}	Spor	w_{11}	w_{12}	w_{13}	Spor	w_{11}	w_{12}	w_{13}	Spor	Spor
w_{21}	w_{22}	w_{23}	Ekonomi	w_{21}	w_{22}	w_{23}	Ekonomi	w_{21}	w_{22}	w_{23}	Ekonomi	
w_{11}	w_{12}	w_{13}	Spor	w_{11}	w_{12}	w_{13}	Spor	w_{11}	w_{12}	w_{13}	Spor	Ekonomi
w_{21}	w_{22}	w_{23}	Ekonomi	w_{21}	w_{22}	w_{23}	Ekonomi	w_{21}	w_{22}	w_{23}	Ekonomi	



w_{k1}	w_{k2}	w_{k3}	?	w_{k1}	w_{k2}	w_{k3}	?	w_{k1}	w_{k2}	w_{k3}	?	?
----------	----------	----------	---	----------	----------	----------	---	----------	----------	----------	---	---

VERİ MADENCİLİĞİ (2)

- K-en yakın komşu
- Destek vektör makineleri
- Naive Bayes
- Yapay Sinir Ağları
- Karar Ağaçları...

DEĞERLENDİRME (1)

- Doğruluk, doğru sınıflandırılan kayıtların sayısının yanlış sınıflandırılan kayıtların sayısına oranı olarak tanımlanmaktadır.
- Doğruluk, diğer bir adıyla sınıflandırıcının doğru tahmin oranıdır. Kesinlik (p); gerçek sınıfı ve tahmin edilen sınıfı 1 olan kayıtların, tahmin edilen sınıfı 1 olan kayıtlara oranı şeklinde tanımlanmaktadır.
- Duyarlılık (r), gerçek sınıfı ve tahmin edilen sınıfı 1 olan kayıtların gerçek sınıfı 1 olan kayıtlara oranıdır.
- F-ölçümü kesinlik ve duyarlılık ölçümlerinin harmonik ortalaması alınarak bulunmektedir.

DEĞERLENDİRME (2)

GERÇEK SINIF	TAHMİN EDİLEN SINIF		
		Sınıf=1	Sınıf=0
	Sınıf=1	a (Doğru Pozitif)	b (Yanlış Negatif)
Sınıf=0	c (Yanlış Pozitif)	d (Doğru Negatif)	

$$\text{Doğruluk} = \frac{a + d}{a + b + c + d}$$

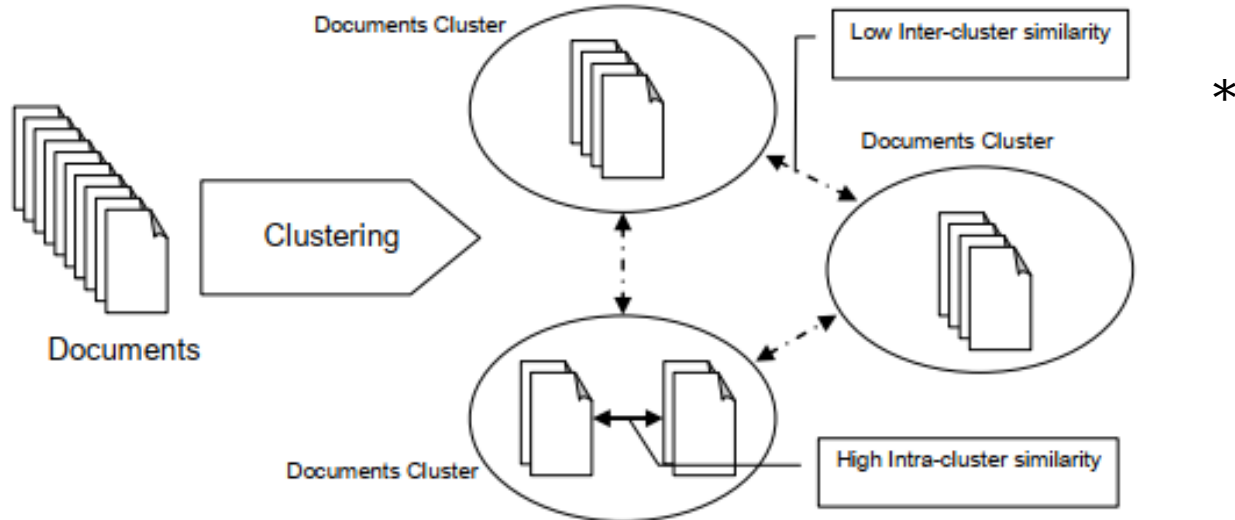
$$\text{Duyarlılık} = \frac{a}{a + b}$$

$$\text{Kesinlik} = \frac{a}{a + c}$$

$$F - \text{ölçümü} = \frac{2 \times \text{Kesinlik} \times \text{Duyarlılık}}{\text{Kesinlik} + \text{Duyarlılık}}$$

VERI MADENCİLİĞİ (3)

- **Kümeleme:** Doküman koleksiyonunda yer alan dokümanları kümeler altında gruplar. Küme içi benzerlik maksimum iken, kümeler arası benzerlik minimum olmalıdır.
- Sınıflandırmanın aksine kümeleme yapılacak veri kümesindeki dokümanlarının sınıf etiketi bulunmamaktadır.



VERİ MADENCİLİĞİ (4)

- K-means
- Hiyerarşik kümeleme,...

DEĞERLENDİRME (3)

- Kümeler içi benzerlik maksimum, kümeler arası benzerlik minimum olması gerekmektedir. Temel değerlendirme ölçütü bu kuraldır.

PYTHON KÜTÜPHANELERİ

- **NLTK (Natural Language Toolkit)** : Önişleme adımlarının gerçekleştirilmesini sağlayan kütüphanedir.
- **Spacy**: NLTK ile aynı görevleri gerçekleştirmektedir.
- **Scikit-learn**: Makine öğrenmesi yöntemlerini sunan kütüphanedir. Ayrıca metin önişleme görevlerinin de yerine getirilmesini sağlamaktadır.
- **Gensim**: Konu modelleri, vektör uzayı modellerini sunan kütüphanedir.