

Bilgiye Eriřim Sistemleri Information Retrieval (IR) Systems

M.Fatih AMASYALI
BLM 5212 Doęal Dil İřlemeye Giriř Ders Notları

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĐİ BÖLÜMÜ



Akış

- Örnek IR Sistemleri
- IR Sistem Mimarisi
- Arama Motoru Mimarisi
- Vektör Uzayı
- Google
- Anahtar Kelime Problemleri
- Zeki IR Teknikleri
- IR Sistemlerinin Deęerlendirilmesi

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĐİ BÖLÜMÜ



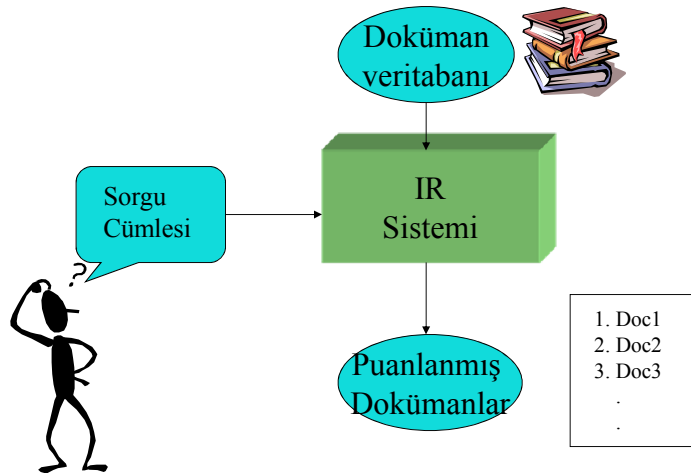
Örnek IR Sistemleri

- Kütüphane veritabanları
anahtar kelime, başlık, yazar, konu vs. ile büyük veritabanlarında arama (www.library.unt.edu)
- Metin Tabanlı Arama Motorları (Google, Yahoo, Altavista vs).
Anahtar kelimelerle arama
- Multimedya Arama (QBIC, WebSeek, SaFe)
Görsel öğelerle arama (şekil, renk vs.)
- Soru Cevaplama Sistemleri (AskJeeves, Answerbus)
Doğal dille arama

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

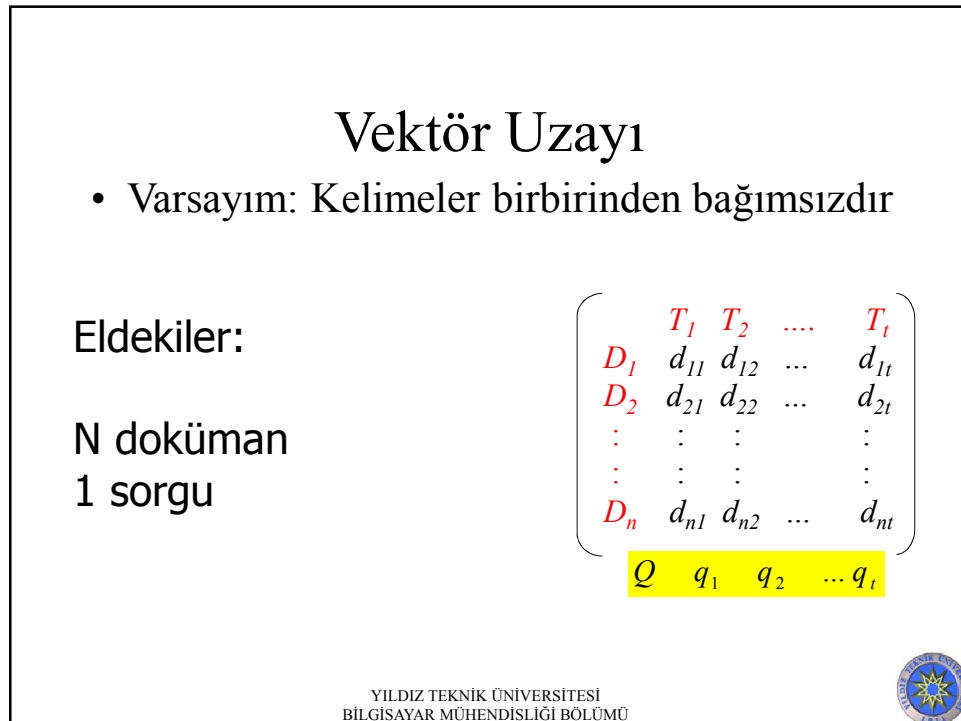
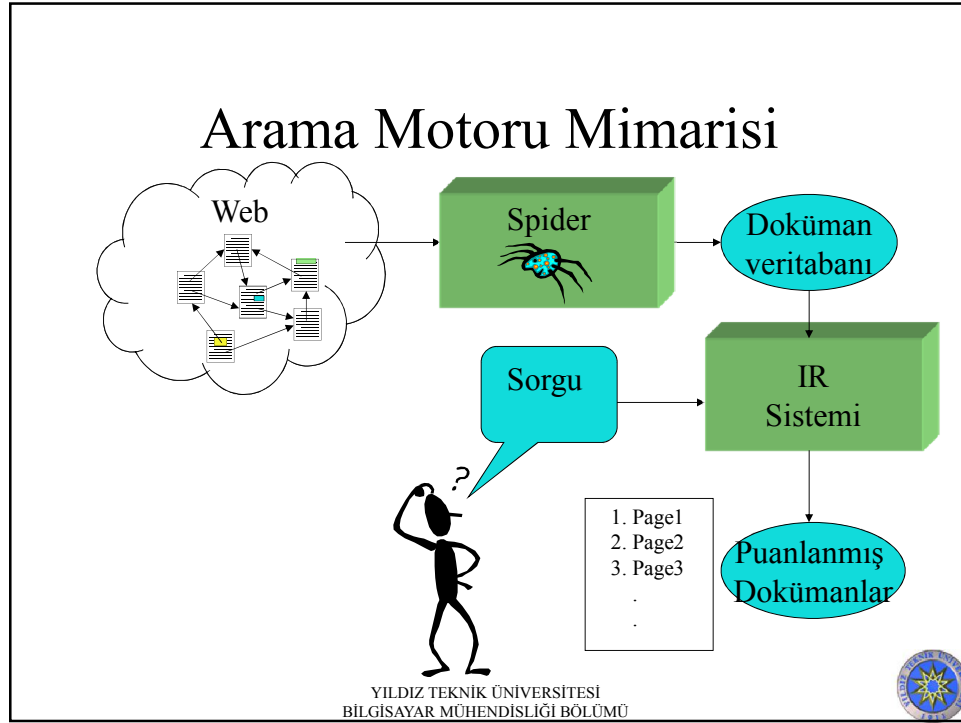


IR Sistem Mimarisi



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ





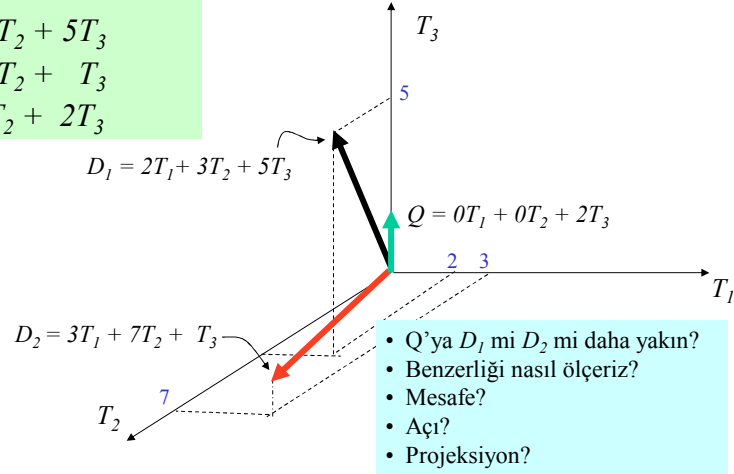
Grafiksel Gösterim

Örnek:

$$D_1 = 2T_1 + 3T_2 + 5T_3$$

$$D_2 = 3T_1 + 7T_2 + T_3$$

$$Q = 0T_1 + 0T_2 + 2T_3$$



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Benzerlik Ölçümü- Inner Product

$$\text{sim}(D_i, Q) = \sum_{k=1}^t (D_i \cdot Q)$$

$$= \sum_{j=1}^t d_{ij} * q_j$$

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Inner Product - Örnek

Binary:

- D = $\begin{matrix} \text{retrieval} & \text{database} & \text{architecture} & \text{computer} & \text{text} & \text{management} & \text{information} \\ \text{1,} & \text{1,} & \text{1,} & \text{0,} & \text{1,} & \text{1,} & \text{0} \end{matrix}$ • Vektör boyutu =
 - Q = $\begin{matrix} \text{1,} & \text{0,} & \text{1,} & \text{0,} & \text{0,} & \text{1,} & \text{1} \end{matrix}$ Sözlük boyutu = 7

→ $\text{sim}(D, Q) = 3$

Ağırlıklı

$$D_1 = 2T_1 + 3T_2 + 5T_3$$

$$Q = 0T_1 + 0T_2 + 2T_3$$

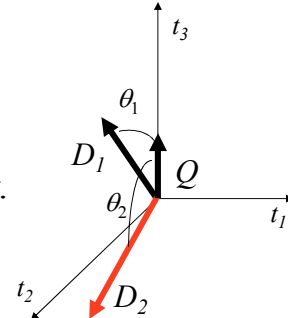
$$\text{sim}(D_1, Q) = 2*0 + 3*0 + 5*2 = 10$$

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Cosine Benzerlik Ölçümü

- İki vektör arasındaki açının cosinüsü
- Inner product, vektör büyüklükleriyle normalize edilir.



$$\text{CosSim}(D_i, Q) = \frac{\sum_{k=1}^t (d_{ik} \cdot q_k)}{\sqrt{\sum_{k=1}^t d_{ik}^2 \cdot \sum_{k=1}^t q_k^2}}$$

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Cosine Benzerlik: Örnek

$$\begin{aligned} D_1 &= 2T_1 + 3T_2 + 5T_3 & \text{CosSim}(D_1, Q) &= 0.81 \\ D_2 &= 3T_1 + 7T_2 + T_3 & \text{CosSim}(D_2, Q) &= 0.13 \\ Q &= 0T_1 + 0T_2 + 2T_3 \end{aligned}$$



Doküman ve Terim Ağırlıkları

- Ağırlıklar dokümanlardaki frekanslarla (tf) ve tüm doküman kütüphanesindeki frekanslarla (idf) hesaplanır.

$tf_{ij} = j$. terimin i . dokümandaki frekansı

$df_j = j$. terimin doküman frekansı

= j . terimi içeren doküman sayısı

$idf_j = j$. terimin ters doküman frekansı

= $\log_2(N/df_j)$ (N : toplam doküman sayısı)



Terim Ağırlıklarının Bulunması

- j . terimin i . doküman için ağırlığı:

$$d_{ij} = tf_{ij} \bullet idf_j = tf_{ij} \bullet \log_2 (N/df_j)$$

- TF → Terim Frekansı

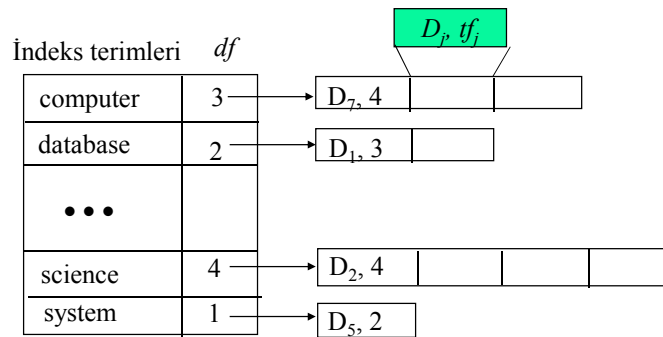
- Bir dokümanda sıkça geçen ancak diğer dokümanlarda pek bulunmayan terimin ağırlığı yüksek olur.
- $\max_i \{tf_{i1}\} = i$. dokümanda en çok geçen terimin frekansı
- Normalizasyon: terim frekansı = $tf_{ij} / \max_i \{tf_{i1}\}$

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Terim Frekansı Uygulaması

- Pratikte doküman vektörleri direkt olarak saklanmaz. Hafıza problemlerinden ötürü, arama için aşağıdaki gibi bir yapıda saklanırlar.



YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Inverted index

- Metinler
- $T_0 = \text{"it is what it is"}$
- $T_1 = \text{"what is it"}$
- $T_2 = \text{"it is a banana"}$

"what", "is" ve "it" kelimeleriyle arama yapılırsa.

$$\{0, 1\} \cap \{0, 1, 2\} \cap \{0, 1, 2\} = \{0, 1\}$$

inverted file index:

- "a": {2}
- "banana": {2}
- "is": {0, 1, 2}
- "it": {0, 1, 2}
- "what": {0, 1}

Full inverted file index: (pozisyonları da içerir)

- "a": {(2, 2)}
- "banana": {(2, 3)}
- "is": {(0, 1), (0, 4), (1, 1), (2, 1)}
- "it": {(0, 0), (0, 3), (1, 2), (2, 0)}
- "what": {(0, 2), (1, 0)}

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ

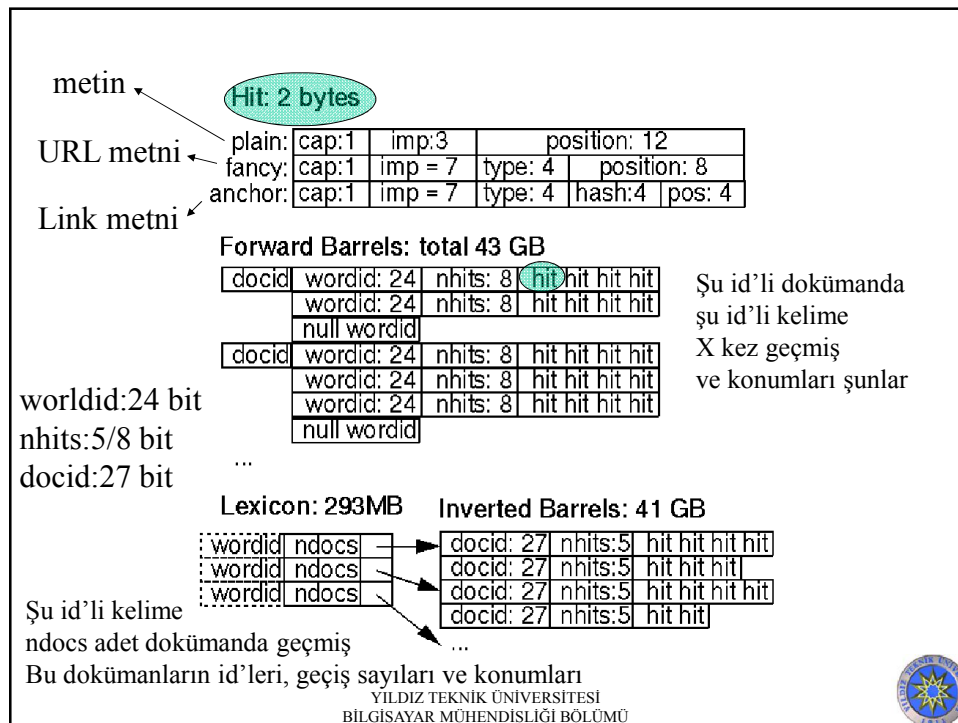
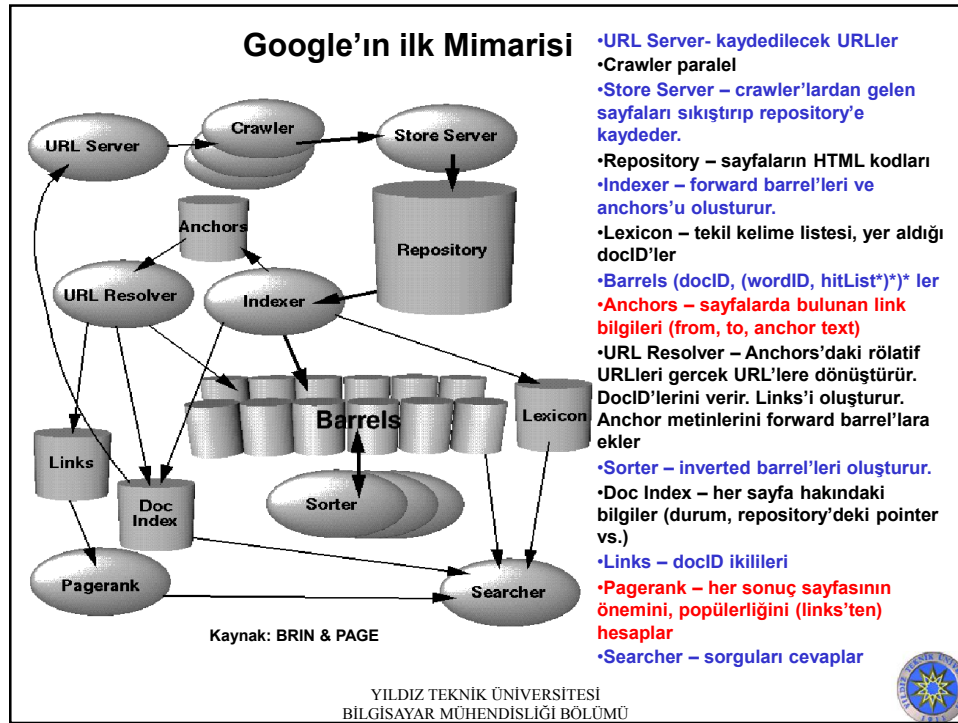


Web Katalogları vs. Arama Motorları

- | | |
|--|---|
| <ul style="list-style-type: none"> • Web Katalogları <ul style="list-style-type: none"> – Elle seçilmiş siteler – Sayfaların içeriğinde değil, tanımlarında arama – Hiyerarşik kategorilere atanırlar | <ul style="list-style-type: none"> • Arama Motorları <ul style="list-style-type: none"> – Tüm sitelerdeki tüm sayfalar – Sayfaların içeriğinde arama – Sorgu geldikten sonra bulunan skorlara göre sıralanırlar. |
|--|---|

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ





Puanlama Sistemi

- Kriterler
 - Pozisyon, Font Büyüklüğü, Büyük Harfle/ Bold/italik yazılma
 - Sitenin popülerliği (PageRank)
 - Başlık ,link metni (Anchor Text), URL metni vs.

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



İlk Google'ın depo büyüklükleri Web istatistikleri

Total Size of Fetched Pages	147.8 GB	Number of Web Pages Fetched	24 million
Compressed Repository	53.5 GB	Number of URLs Seen	76.5 million
Short Inverted Index	4.1 GB	Number of Email Addresses	1.7 million
Full Inverted Index	37.2 GB	Number of 404's	1.6 million
Lexicon	293 MB		
Temporary Anchor Data (not in total)	6.6 GB		
Document Index Incl. Variable Width Data	9.7 GB		
Links Database	3.9 GB		
Total Without Repository	55.2 GB		
Total With Repository	108.7 GB		

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Sistem Performansı

- 26 milyon site 9 günde indirilmiş.(Saniyede 48.5 sayfa)
- Indexer ve Crawler aynı anda çalışıyor
- Indexer saniyede 54 sayfayı indeksliyor
- Sorter'lar 4 makinede paralel çalışarak 24 saatte inverted index'i oluşturuyor



Anahtar Kelime ile Aramada Problemler

- Eş anlamlı kelimeleri içeren dokümanlar bulunamaz.
 - “PRC” vs. “China”
- Eşsesli kelimeler ilgisiz dokümanların bulunmasına sebep olabilir.
 - “bat” (baseball vs. mammal)
 - “Apple” (company vs. fruit)
 - “bit” (unit of data vs. act of eating)



Zeki IR Teknikleri

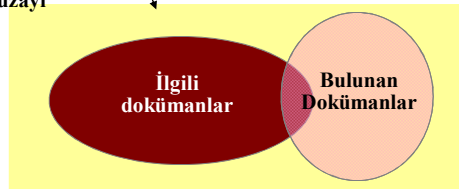
- Kelimelerin anlamları
- Sorgudaki kelimelerin sırası
- Kullanıcılardan döndürülen sonuçların kalitesiyle ilgili alınan geri bildirimler (sonuçların kaçınıcısına tıkladı, kaç sonuç sayfası inceledi vb.)
- Aramayı ilgili kelimelerle genişletmek
- İmla denetimi – kelime önermek
- Kaynakların güvenilirliği
- Kişiselleştirilmiş arama
- Eklemeli dillerde bazı eklerden bağımsızlık

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Değerlendirme

Tüm doküman
uzayı



relevant	retrieved & irrelevant	Not retrieved & irrelevant
	retrieved & relevant	not retrieved but relevant
	retrieved	not retrieved

$$\text{recall} = \frac{\text{Number of relevant documents retrieved}}{\text{Total number of relevant documents}}$$

$$\text{precision} = \frac{\text{Number of relevant documents retrieved}}{\text{total Number of documents retrieved}}$$

YILDIZ TEKNİK ÜNİVERSİTESİ
BİLGİSAYAR MÜHENDİSLİĞİ BÖLÜMÜ



Kaynaklar

- <http://www.cs.huji.ac.il/~sdbi/2000/google/index.htm>
- Searching the Web ,Ray Larson & Warren Sack
- Knowledge Management with Documents, Qiang Yang
- Introduction to Information Retrieval, Rada Mihalcea
- Wikipedia
- Introduction to Information Retrieval, Evren Ermis
- The Anatomy of a Large-Scale Hypertextual Web Search Engine, Sergey Brin, Lawrence Page (<http://infolab.stanford.edu/~backrub/google.html>)

