

Hesaplamalı Anlambilim Ders Notları

# Anlamsal Uzaylar

Mehmet Fatih AMASYALI  
Hesaplamalı Anlambilim  
Ders Notları  
BÖLÜM 3



Yıldız Teknik Üniversitesi  
Bilgisayar Mühendisliği Bölümü

NOVA Research Lab

Hesaplamalı Anlambilim Ders Notları

## İçerik

- İstatistiksel Dil Modelleri
- Kelimelerin Anlamsal Benzerliklerinin Ölçümü
  - Gizli Anlam İndeksleme (Latent Semantic Indexing - LSI)
  - Latent Dirichlet Allocation - LDA
  - Kavramsal hiyerarşilerin kullanımı
  - Birbiri yerine kullanılan kelimeler
  - Kelime kümeleme
  - Anlamsal Uzaylar
  - Kelime öbekleri



NOVA Research Lab

## İstatistiksel Dil Modelleri\*

- Bir cümlenin olasılığını bulmak
- Kullanım alanları
  - Makine Çevirisi
    - $P(\text{taze balık aldım}) > P(\text{yeni balık aldım})$
  - Yazım düzeltimi
    - Ali soan yedi
    - $P(\text{ali soğan yedi}) > P(\text{ali sokaan yedi})$
  - Konuşmadan Metne
    - $P(\text{soğan yedim}) > P(\text{sol an yedim})$
  - vb.



[https://web.stanford.edu/~jurafsky/slp3/slides/LM\\_4.pdf](https://web.stanford.edu/~jurafsky/slp3/slides/LM_4.pdf)

NOVA Research Lab

## İstatistiksel Dil Modelleri

- Amaç 1: bir cümlenin olasılığını bulmak  
 $P(W) = P(w_1, w_2, \dots, w_n)$
- Amaç 2: bir kelime sekansından sonra gelecek kelimenin olasılığını bulmak  
 $P(w_5 | w_1, w_2, w_3, w_4)$   
 Bunları nasıl hesaplayacağız?  
 Bayes Kuralı :  $P(A|B) = P(A, B) / P(B) \rightarrow$   
 $P(A, B) = P(B) * P(A|B) = P(A) * P(B|A)$   
 Zincir Kuralı:  
 $P(A, B, C, D) = P(A) * P(B|A) * P(C|A, B) * P(D|A, B, C)$



NOVA Research Lab

## İstatistiksel Dil Modelleri

- Zincir kuralının genel hali:
- $P(w_1, w_2, w_3, \dots, w_n) =$   
 $P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \dots P(w_n|w_1, \dots, w_{n-1})$   
 Örnek:  $P(\text{ali}, \text{okuldan}, \text{buraya}, \text{geldi}) = P(\text{ali}) * P(\text{okuldan}|\text{ali}) * P(\text{buraya}|\text{ali}, \text{okuldan}) * P(\text{geldi}|\text{ali}, \text{okuldan}, \text{buraya})$
- $P(\text{geldi}|\text{ali}, \text{okuldan}, \text{buraya})$  yı nasıl bulacağız?

Bir olasılık :

$$P(\text{geldi}|\text{ali}, \text{okuldan}, \text{buraya}) = \frac{\text{frekans}(\text{ali okuldan buraya geldi})}{\text{frekans}(\text{ali okuldan buraya})}$$

Ama bu yolla güvenilir değerler elde etmek için yeterli büyüklükte bir derlemimiz elimizde olmaz.

Çoğu frekans 0 olacaktır. Çünkü dil, diskrit uzayda çok seyrek. Peki ne yapalım?



## İstatistiksel Dil Modelleri

- Markov imdada yetişir ☺
- Markov varsayımı: her olasılık kendinden önceki **k** bileşene bağlıdır.
- **k=0 için** her olasılık bağımsızdır
  - $P(w_1, w_2, w_3, w_4) = P(w_1) * P(w_2) * P(w_3) * P(w_4)$
  - $P(w_4|w_1, w_2, w_3) = P(w_4)$
  - $P(\text{geldi}|\text{ali}, \text{okuldan}, \text{buraya}) = P(\text{geldi})$
- **k=1 için** her olasılık sadece bir öncesine bağlıdır
  - $P(w_1, w_2, w_3, w_4) = P(w_2|w_1) * P(w_3|w_2) * P(w_4|w_3)$
  - $P(w_4|w_1, w_2, w_3) = P(w_4|w_3)$
  - $P(\text{geldi}|\text{ali}, \text{okuldan}, \text{buraya}) = P(\text{geldi}|\text{buraya})$



- **k=2 için** her olasılık sadece bir ve iki öncekine bağlıdır
  - $P(w_1, w_2, w_3, w_4) = P(w_3|w_1, w_2) * P(w_4|w_2, w_3)$
  - $P(w_4|w_1, w_2, w_3) = P(w_4|w_2, w_3)$
  - $P(\text{geldi}|\text{ali}, \text{okuldan}, \text{buraya}) = P(\text{geldi}|\text{okuldan}, \text{buraya})$
- k değerini istediğimiz gibi arttırabiliriz.



- Markov varsayımı ne yapıyor?
- Varsayım yokken
- $P(w_1, w_2, w_3, \dots, w_n) = P(w_1)P(w_2|w_1)P(w_3|w_1, w_2) \dots P(w_n|w_1, \dots, w_{n-1})$
- $P(w_1, w_2, \dots, w_n) = \prod_{i=1}^n P(w_i|w_1, w_2, \dots, w_{i-1})$
- Varsayım varken
- $P(w_i|w_1, w_2, \dots, w_{i-1}) = P(w_i|w_1, w_2, \dots, w_{i-k-1}, w_{i-k}, w_{i-k+1}, \dots, w_{i-1}) \approx P(w_i|w_{i-k}, w_{i-k+1}, \dots, w_{i-1})$



kırmızıları yok sayıyor

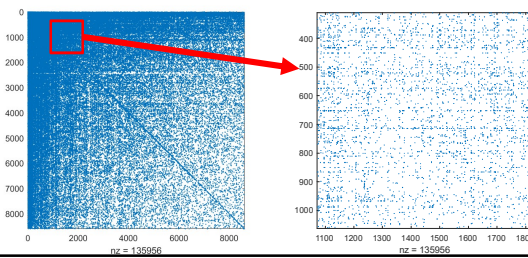
- Markov varsayımının avantajları:
  - daha az işlem
  - daha az seyrek veriler
- Dezavantajı:
  - Dil, uzak bağımlılıklar içerdiğinden çok yeterli bir model değil
    - **Çine** yaptığım ziyarette ... çok uğraştım ama bir türlü **Çinceyi** öğrenemedim.
    - Burada **Çinceyi** kelimesi belki onlarca kelime öncesindeki **Çine** kelimesine bağlı.



K seçimi burada bir ikilem

NOVA Research Lab

- **Kod dil\_modeli.m**
- 1000 ekonomi haberi
- 5 ten az geçen kelimeleri sildikten sonra yaklaşık 8K\*8K lık bir matris
- Çok seyrek, ~64M den sadece ~140K 0 değil



NOVA Research Lab

## Hesaplamalı Anlambilim Ders Notları

- **ekonomi**'den sonra en fazla geçen 5 kelime:  
 $P(w_i|ekonomi)$  değeri en yüksek 5 kelime:  
bakanı, ve, ligi, bakanlığı, bankası
- **para**'dan sonra: politikası, cezası, ve, birimlerinin
- **çok**'dan sonra: önemli, daha, büyük, sayıda, iyi



NOVA Research Lab

## Hesaplamalı Anlambilim Ders Notları

- Bir cümleinin bigram modeline göre olasılığı:
- $P(\text{fuvar otomotiv sektörüne önemli bir hareket getirdi}) = P(\text{otomotiv}|\text{fuvar}) * P(\text{sektörüne}|\text{otomotiv}) * P(\text{önemli}|\text{sektörüne}) * P(\text{bir}|\text{önemli}) * P(\text{hareket}|\text{bir}) * P(\text{getirdi}|\text{hareket})$
- $P(w_i|w_{i-1}) = \frac{fr(w_{i-1}, w_i)}{fr(w_{i-1})}$
- $P(\text{otomotiv}|\text{fuvar}) = \text{fre}(\text{fuvar otomotiv}) / \text{fre}(\text{fuvar})$
- $\text{Log}(P(\text{fuvar otomotiv sektörüne önemli bir hareket getirdi})) = -21.2067$

$w_i$	$fr(w_{i-1})$	$fr(w_{i-1}, w_i)$	$P(w_i w_{i-1})$
fuvar	18	1	1/18
otomotiv	78	4	4/78
sektörüne	35	1	1/35
önemli	560	162	162/560
bir	3940	4	4/3940
hareket	42	1	1/42



NOVA Research Lab

## 0'larla başlamak Smoothing

- $\text{Log}(P(\text{fuar otomotiv sektörüne önemli bir hareket getirdi})) = \text{log}(P(\text{otomotiv|fuar})) + \text{log}(P(\text{sektörüne|otomotiv})) + \text{log}(P(\text{önemli|sektörüne})) + \text{log}(P(\text{bir|önemli})) + \text{log}(P(\text{hareket|bir})) + \text{log}(P(\text{getirdi|hareket})) = -21.2067$
- $\text{Log}(P(\text{fuar otomotiv sektörüne önemli iki hareket getirdi})) = -\text{INF}$
- **Add one estimation = Laplace smoothing**
- $P_{LS}(w_i|w_{i-1}) = \frac{fr(w_{i-1}, w_i) + 1}{fr(w_{i-1}) + V}$
- $\text{Log}(P_{LS}(\text{fuar otomotiv sektörüne önemli bir hareket getirdi})) = -44$
- $\text{Log}(P_{LS}(\text{fuar otomotiv sektörüne önemli iki hareket getirdi})) = -49$



## En olası dizilişi bulmak

- Kelimeler = önemli, hareket, bir
- Olası dizilişler
- $P_{LS}(\text{önemli bir hareket}) = -11.8570$
- $P_{LS}(\text{önemli hareket bir}) = -18.1881$
- $P_{LS}(\text{bir önemli hareket}) = -17.4616$
- $P_{LS}(\text{bir hareket önemli}) = -16.8926$
- $P_{LS}(\text{hareket önemli bir}) = -13.0944$
- $P_{LS}(\text{hareket bir önemli}) = -17.4034$
- Smoothing kullanmasaydık?
- Bir başka uygulama: Şu kelimelerle bir cümle kurun



## Dil modellerini deęerlendirmek

- Elimizde 2 dil modeli olsun (A, B). Hangisi daha iyi?
- Harici yöntem: her 2 modeli de farklı görevler için (Makine Çevirisi, Yazım düzeltimi, Konuşmadan Metne) kullan. O görevlerdeki performanslarına göre karşılaştır
- Dahili yöntem: Derlemi eğitim ve test kümesi olarak ayır. A ve B yi eğitim üzerinden oluştur.
  - Test kümesindeki cümlelere hangisi daha yüksek olasılık veriyorsa o daha iyidir.
  - Bir cümle başı verildiğinde sonunu hangisi doğru tahmin ediyorsa o daha iyidir.
  - Bu görevlerde Markov varsayımında  $k=0$  seçimi nasıldır?



## Dil modelleriyle cümle / dizilim üretimi

- Bir ilk kelime seçelim. Bundan sonra gelecek kelimeleri dil modeli ile belirleyelim.
- $K=0$  için her zaman en yüksek olasılığı seçersek hep aynı kelime tekrar eder
- $K=1$  için, her zaman en yüksek olasılığı seçersek bir kelimeden sonra hep aynı kelime gelir.
- $K=2$  için, ardışık 2 kelimeden sonra hep aynı 3. kelime gelir.
- Bunu aşmanın 2 yolu var
  - En yüksek olasılıklı  $t$  taneden birini rasgele seçmek
  - Hepsini içinden olasılıklarına göre seçim yapmak (Shannon Game)
- Ne zaman duracağız: İsteddiğimiz sayıda kelime / cümle ürettiğimizde





## Hesaplamalı Anlambilim Ders Notları

- **Kod dil\_modeli.m**
- Bigram ( $K=1$ ) modeli ile
- En yüksek olasılıklı  $t$  taneden birini rasgele seçer.  $t$  değeri azaldıkça metinlerde kopya çekme olasılığı artar.  $t=5$
- 1000 ekonomi haberi ile
  - "fonların yüzde 5 de bu konuda son 10 ın da çok sayıda türk telekom gibi çok önemli olduğunu belirten bakan şimşek 2012 nin."
  - "mesleki eğitim ve ticaret bakanlığı ın piyasa ile ilgili olarak belirlendi bu nedenle bir önceki yıla da en fazla artış oranı ise lira oldu bu."
  - "akademik araştırmalar sonucunda yer aldığı bilgiye göre bir şekilde diyen ve yüzde 10 yılda bir süre içerisinde türkiye istatistik verilerine borsa seçimlerinin dikkati çekti türkiye."
- Burada first order ( $K=1$ ) yapının problemi bariz görünüyor. "yer aldığı bilgiye" "yüzde 10 yılda"



NOVA Research Lab

## Hesaplamalı Anlambilim Ders Notları

- **Kod dil\_modeli\_second.m**
- trigram ( $K=2$ ) modeli,  $t=5$
- 1000 ekonomi haberiyle üretilen cümleler
  - "yaşanan sıkıntıları için devamlı ucuz bilet satıyoruz 100 90 ı aslında biz bu bitireceğiz inşallah yıl sonuna kadar türkiyede üretilen fazlasını belirterek geçen yıl kasım ayı ihracatını"
  - "birlikte fiyatlardaki üzerindeki etkileri takip edilecek gerek halinde düzenlemede birtakım yeni getirilecek talepler daha önceki bir açıklamasında 2012 yılının ilk 2 ihaleyi 500 adet olacak satış faaliyet"
- 1000 magazin haberiyle üretilen cümleler:
  - "arkadaşlarının da kendisini hiç söyleyen nilüfer kanser olduğunu tüm anlattı ilk itibaren dizinin kerem benim olduğu gibi bu hafta muhteşem neler olacak ali ikna eder ama şimdi"
  - "oyuncuların eğitimi yoktu ama hep başka şeyleri yapmak için burada bu arkadaşlar benim dönüş ve benim gibi tabii ki oldu ama devam çok bir adam alman sevgiliyle"
- Ardışık 3 lüler daha anlamlı artık.



NOVA Research Lab

- [Kod dil\\_modeli\\_term\\_doc.m](#)
- Olasılık hesabında ardışık geçme yerine **aynı dokümanda geçmeyi** kullansak?
- Frekansı çok yüksek kelimeleri (stopWord) atmak gerekli. Yoksa her kelimeye en benzer kelimeler onlar olur.
- $K=1$  ve  $t=15$  için üretilen cümle: "faiz oranı yüksek yönetim kurumsal yönetim kurulu genel 2 satış 5 enflasyon merkez cari kredi yapı nın türk türkiyenin ancak cari büyüme merkez bankası merkez cari."
- Cümle üretimi için iyi değil. Birlikte geçen kelimeleri sıralıyor.
- Faiz e en benzeyen kelimeler : "merkez" "kredi" "bankası" "oranı" "cari"
- Benzer anlama sahip kelimeleri bulmak için kullanılabilir. İleride göreceğiz.



## Kelime benzerliği

- Karakter benzerliği (edit-distance)
  - Levenshtein vb. [kod lev.m](#), [lev\\_sample.m](#)
- Anlam benzerliği
  - (Cimbom, Galatasaray,?)

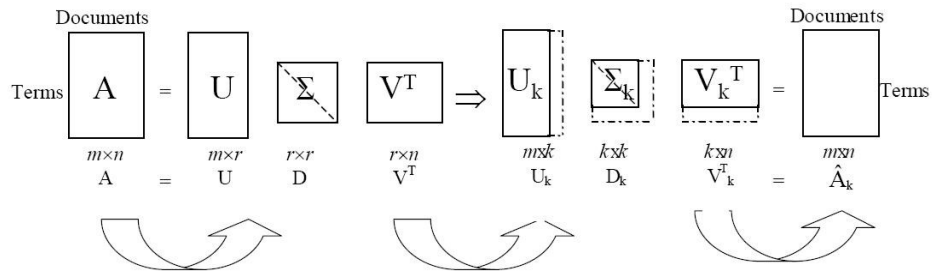


Hesaplamalı Anlambilim Ders Notları

# Saklı Anlam İndeksleme

## Latent Semantic Indexing

- Kelimelerin ardındaki saklı kavramları ortaya çıkarıp dokümanları bu kavramlar uzayında temsil etmek



\* <https://liqiangguo.wordpress.com/2011/06/09/latent-semantic-analysis/>

NOVA Research Lab

Hesaplamalı Anlambilim Ders Notları

# Saklı Anlam İndeksleme

```
% 2 sinifa ait dokumanlarımız olsun.
```

```
% 1. sınıf kedi 2. sınıf köpek
```

```
% kedi miyav süt fare köpek havla kemik hayvan tüy
```

```
D=[ 1 1 0 0 0 0 0 0 0 ; % kedi sinifi
    1 0 1 0 0 0 0 0 0 ; % kedi sinifi
    0 0 1 0 0 0 0 1 0 ; % kedi sinifi
    1 0 0 0 0 0 0 0 1 ; % kedi sinifi
    0 0 0 0 1 1 0 0 1 ; % köpek sinifi
    0 0 0 0 1 0 0 1 0 ; % köpek sinifi
    0 0 0 0 0 1 1 0 0 ; % köpek sinifi
    0 0 0 0 1 0 1 1 0 ]; % köpek sinifi
```

```
D=D'; % matrisin kelime*dokuman olması gerek
```

```
[U,S,V] = svd(D);
```

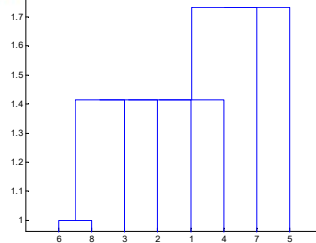
```
%D = U*S*V'
```

```
%D 9*8 kelime*dokuman
```

```
%U 9*9 kelime*kavram
```

```
%S 9*8 singular values
```

```
%V 8*8 dokuman*kavram
```



NOVA Research Lab

## Saklı Anlam İndeksleme

```
r=5; % dokumanlari kac gizli kavram boyutuna indirgeycegimiz
```

```
V=V';
```

```
yeniD= U(:,1:r)*S(1:r,1:r)*V(1:r,:);
```

```
yeniD=yeniD(1:r,:);
```

```
figure;
```

```
Z = linkage(yeniD');
```

```
dendrogram(Z);
```

```
diag(S)
```

```
ans =
```

```
2.4893
```

```
2.0216
```

```
1.8073
```

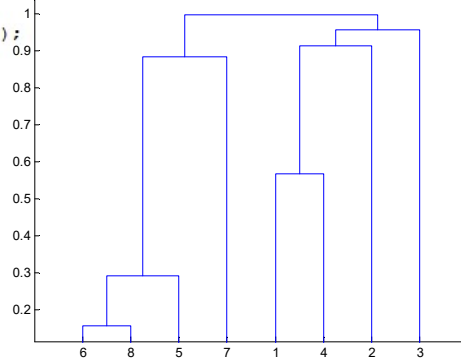
```
1.3486
```

```
1.2349
```

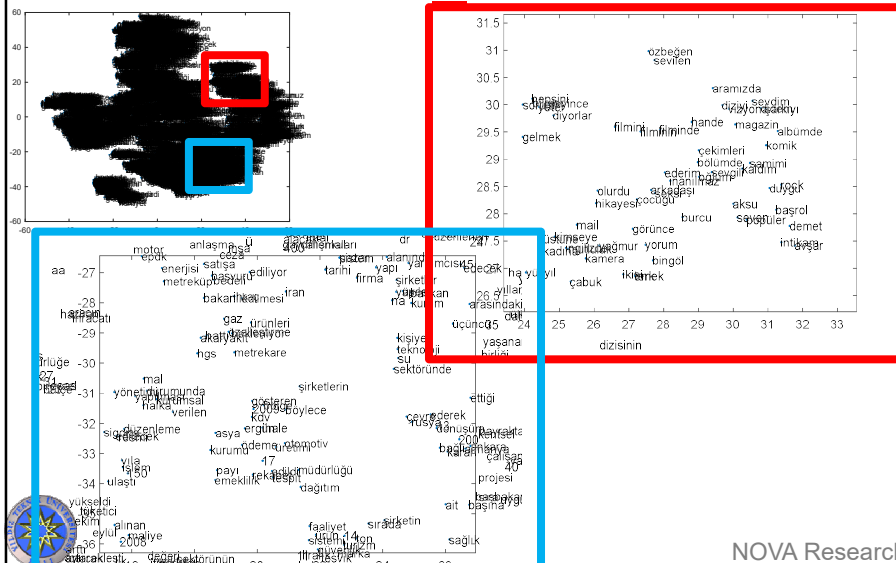
```
0.8247
```

```
0.5943
```

```
0.2694
```



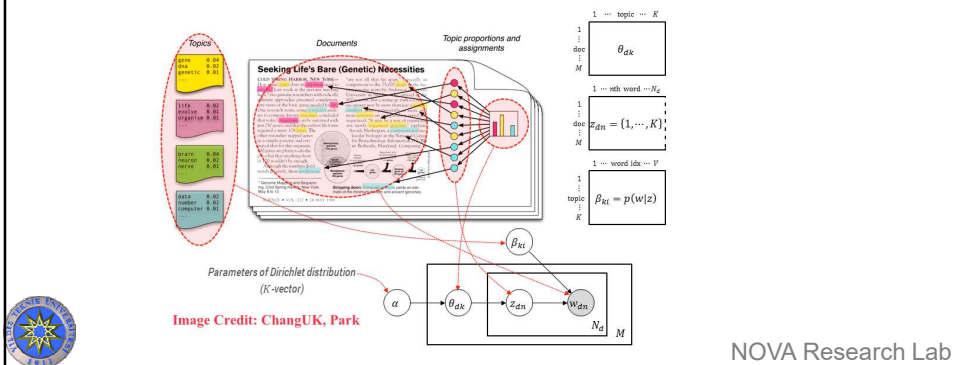
## Büyük ölçekte LSI , kod lsi modeli.m



Hesaplamalı Anlambilim Ders Notları

# Latent Dirichlet Allocation (LDA) model, Topic Modelling

- Konular kelimelerin karışımı
- Dokümanlar konuların karışımı
- <http://littlesaiph.blogspot.com/2012/07/laymans-explanation-of-online-lda.html>



Hesaplamalı Anlambilim Ders Notları

## kod lda\_modeli.m



## Kavramsal Ağaç (hiyerarşi) ile Anlamsal Benzerlik Bulma

- İki temel yaklaşım
  - Bağ sayma
    - Taksonomi (kavramsal ağaç) yeterli
  - Ortak/ Müşterek bilgi (Mutual Information)
    - Taksonomi ve derlem(corpus) gerekli



## Leacock & Chodorow (1998)

$$sim_{LC}(c_1, c_2) = -\log \frac{len(c_1, c_2)}{2L}$$

- ***len(c1,c2)*** iki synset arasındaki en kısa yolun uzunluğu. (*benzerlik değeriyle ters orantılı*)
- ***L***, tüm taksonominin derinliği

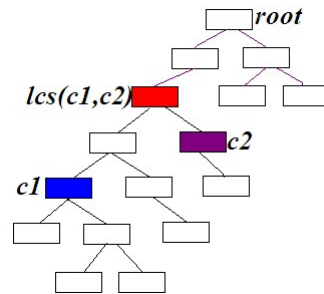


## Wu & Palmer (1994)

$$\text{sim}_{\text{Wu\&Palmer}}(c_1, c_2) = \frac{2 \times N_3}{N_1 + N_2 + 2 \times N_3}$$

• **N1** ve **N2**: **c1** ve **c2**'nin en yakın ortak üst synset'lerine ( **lcs(c1,c2)** ) IS-A bağlarıyla uzaklıkları (benzerlik değeriyle ters orantılı)

• **N3**, en yakın ortak üst synset'in kök synset'e IS-A bağlarıyla uzaklığı (büyüklüğü ortak üst synset'in spesifikliğini gösterir, benzerlikle doğru orantılı)



## Jiang-Conrath (1997)- Lin (1998)

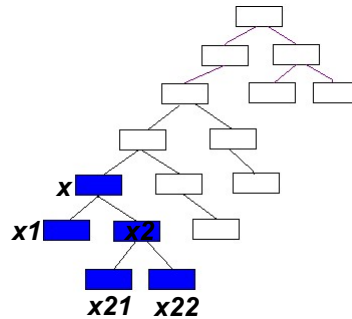
Bir kavramın bir korpus'taki olasılığı =  
corpus'ta geçme sayısı (frekansı)

$$P(\text{concept}) = \text{freq}(\text{concept}) / \text{freq}(\text{root})$$

**freq(x)** → korpusa **x** synset'inin tepesinde bulunduğu

tüm synset'lerin frekanslarının toplamı

$$\text{freq}(x) = f(x) + f(x1) + f(x2) + f(x21) + f(x22)$$



## Jiang-Conrath (1997)- Lin (1998)

$$sim_{JC}(c_1, c_2) = \frac{2 \log(p(lcs(c_1, c_2)))}{\log(p(c_1)) + \log(p(c_2))}$$

$$sim_L(c_1, c_2) = \frac{2 * \log(p(lcs(c_1, c_2)))}{\log(p(c_1)) + \log(p(c_2))}$$

*lcs(c1,c2)* en yakın ortak üst synset

- A, her iki kavramı da içeren en spesifik kavramı kullanır (iki kavramın beraber geçtiği doküman sayısına benzer)
- B, iki kavramdan herhangi birini içeren doküman sayısına benzer



## Karşılaştırma

- Bütün metotlar İngilizce 38 kelime çiftine uygulanmış.
- Bulunan benzerlik değerlerinin, insan yargılarıyla olan korelasyonları yandaki tabloda.

Method	Type	Correlation
Wu & Palmer 1994	Edge Counting	0.74
Li 2003	Edge Counting	0.82
Leacock & Chodorow 1998	Edge Counting	0.82
Resnik 1999	Info. Content	0.79
Lin 1998	Info. Content	0.82
Lord 2003	Info. Content	0.79
Jiang & Conrath 1998	Info. Content	0.83
Tversky 1977	Feature Based	0.73
Adapted Lesk 2002	Feature Based	0.37*
Rodriguez 2003	Hybrid	0.71





Hesaplamalı Anlambilim Ders Notları

## Birbiri yerine kullanılan kelimeler Deniz Yüret, 2009

- W1 W2 X1 W3 W4
- W1 W2 X2 W3 W4
- X1 ve X2 birbiri yerine kullanılmış, aynı bağlamda kullanılmış, demek ki aralarında bir ilişki var.
- X1 ve X2 kaç kere aynı bağlamda kullanıldı ise aralarındaki bağ o kadar güçlü.



NOVA Research Lab

Hesaplamalı Anlambilim Ders Notları

## Birbiri yerine kullanılan kelimeler

İnci Düzenli'nin Bitirme Projesi'nde bulduğu örnekler:

Format: kelime1 kelime2 kaç farklı bağlamda birbirinin yerine kullanıldığı

dolar ytl 939	meclis tbmm 219
bulun olmak 619	belir kaydet 194
avro euro 459	futbol oyun 189
kanun yasa 350	art azal 187
avro dolar 344	avro ytl 182
oyna yap 309	alan saha 177
dolar lira 251	gerek iste 176
dolar euro 249	lira ytl 173
milyar milyon 241	haziran nisan 161



NOVA Research Lab

## Kelime kümeleme

- Kelimeleri
    - terim\*doküman matrisine göre
    - Birlikte geçme matrisine göre
    - Sınıflarda yer alma olasılıklarına göre
- kümelersek aynı kümede yer alan kelimeler birbirine benzerdir diyebiliriz.



## Anlamsal Uzay

- Kelimelerin anlamsal benzerliklerinin bir temsili
- Kavramların anlamsal yakınlıklarına göre yer aldıkları bir uzay
- Anlamsal yakınlığa bağlı koordinatlar nasıl oluşturulur?
  - Terim doküman matrislerini ayrıştırma (LSI)
  - Konu modelleri (LDA)
  - Birlikte geçme matrisleri (MDS, Glove)
  - Sinirsel dil modelleri (word2vec, fasttext)



## Birlikte Geçme Matrisleri

- Harris der ki:
  - Birlikte kullanılan kavramlar birbirlerine anlamsal olarak benzerler.
- Birlikte kullanım: aynı cümlede, aynı metinde, sabit bir kelime penceresinde



## Örnek çalışma-1 Kelimelerin Anlamsal Benzerlik Ölçümü, Amasyalı, 2006

- İki kelimenin Internet'te yer alan sayfaların kaçında yan yana kullanıldıkları bulunarak belirlenmiştir.
- Bunun için arama motoruna  
“kelime1 kelime2” sorgusu gönderilerek gelen sonuç sayfasındaki sonuç sayısı alınmıştır.

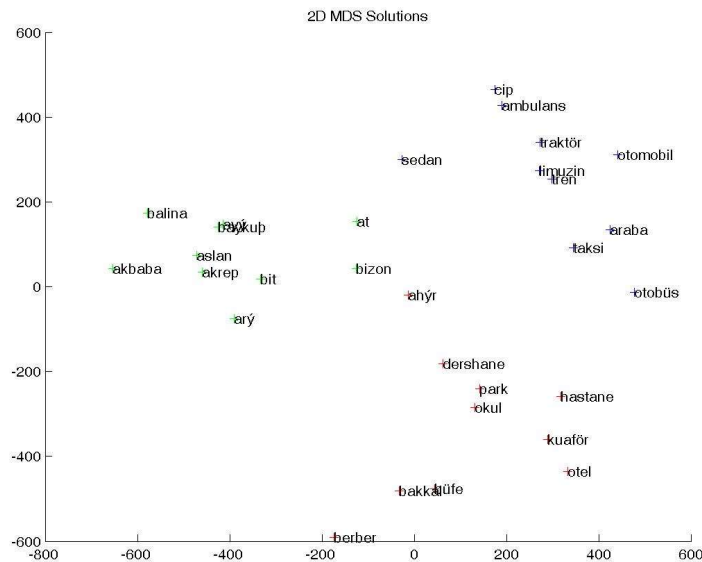


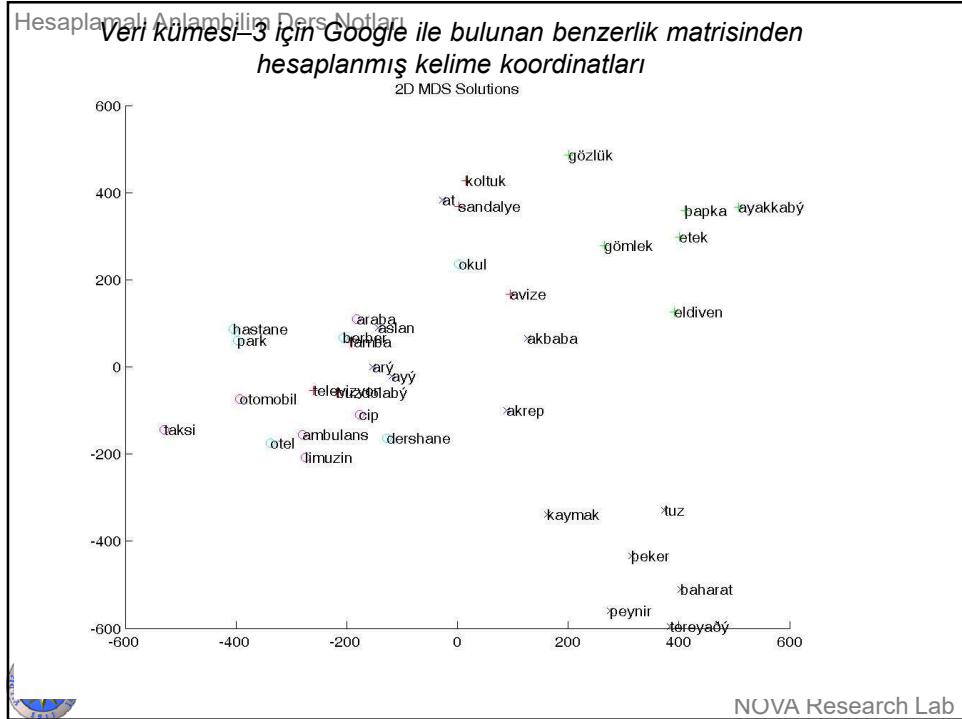
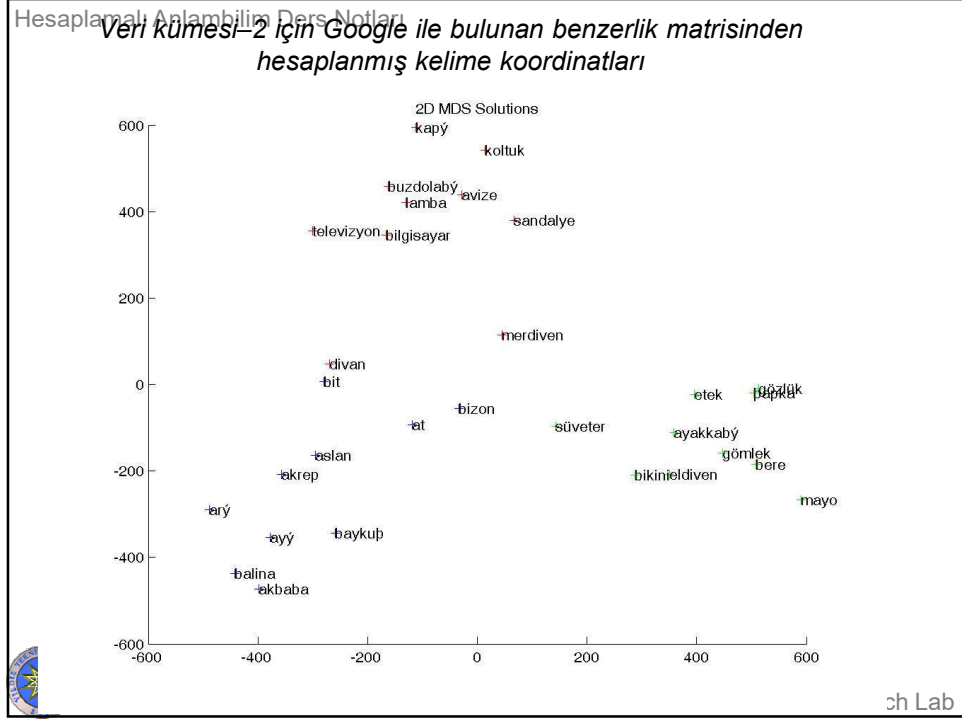
## Kelimeler

Veri kümesi -I	Yer berber – bakkal- kuaför- hastane- otel- okul- dershane- büfe- ahır- park		Hayvan akbaba- arı- baykuş- balina- aslan- at – bizon – akrep - ayı- bit		Taşıt otobüs- cip- taksi- ambulans- araba- sedan- tren- limuzin- otomobil- traktör	
Veri kümesi -II	Ev eşyası divan - televizyon- avize- merdiven- bilgisayar- buzdolabı- kapı- sandalye- lamba- koltuk		Hayvan akbaba- arı- baykuş- balina- aslan- at – bizon – akrep - ayı- bit		Çiyecek bikini- mayo- bere- süveter- etek- ayakkabı- eldiven- gömlek- şapka- gözlük	
Veri kümesi -III	Yiyecek tereyağı kaymak baharat peynir şeker tuz	Çiyecek ayakkabı eldiven şapka gömlek gözlük etek	Ev eşyası avize sandalye lamba koltuk buzdolabı televizyon	Hayvan aslan arı at akrep ayı akbaba	Yer otel dershane park berber hastane okul	Taşıt taksi limuzin ambulans cip otomobil araba



### Veri kümesi-1 için Google ile bulunan benzerlik matrisinden hesaplanmış kelime koordinatları





Hesaplamalı Anlambilim Ders Notları

## Kelimeleri Sınıflandırma Başarıları

	v1	v1	v2	v2	v3	v3
	10 boyut	2 boyut	10 boyut	2 boyut	10 boyut	4 boyut
SVM	96,6	100	83,3	96,6	88,8	77,7
C4.5	83,3	83,3	90	90	77,7	75
RF	90	90	93,3	93,3	72,2	75
EM	66,6	96,6	66,6	96,6	66,6	83,3



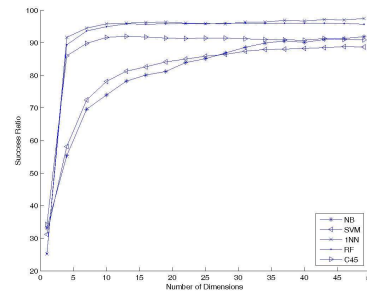
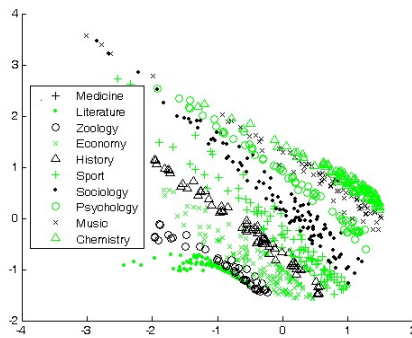
NOVA Research Lab

Hesaplamalı Anlambilim Ders Notları

## Örnek çalışma-2

### Teknik Terim sınıflandırma

- 967 teknik terimin koordinatları, **15 bin web sitesi (neden Google değil?)** üzerinde ile



NOVA Research Lab

### Örnek çalışma-3

#### Metin sınıflandırma (Amasyalı, Beken, 2009)

- 5 farklı haber sınıfına (ekonomi, magazin, sağlık, siyasi, spor) ait 230'ar metin
- Her sınıftan 150'şer haber metni eğitim, 80'er adedi test
- PCKimmo ile kelime kökleri
- farklı gövde sayısı yaklaşık 4500



### Anlamsal Uzayların Metin Sınıflandırmada Kullanımı

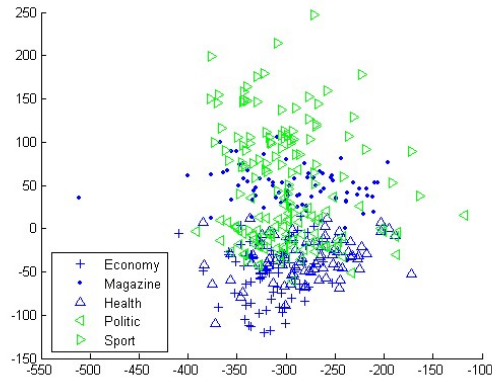
- 4500 kelimenin sayısal karşılıklarını elde etmek için 15.000 web sitesinden oluşan bir külliyat kullanılmış.
- Kelimelerin birbirlerine anlamsal yakınlık matrisi bu külliyatta birlikte geçtikleri doküman sayıları
- uzaklık=1/yakınlık
- MDS (ÇBÖ)
- MDS'de ilk 100 ve ilk 10 boyut alınmış



Hesaplamalı Anlambilim Ders Notları

# Anlamsal Uzayların Metin Sınıflandırmada Kullanımı

- Metinler, içerdikleri kelimelerin koordinatlarının ortalamaları
- Metinler kelimelerle aynı boyutlu
- 1 ve 2. boyutlar



ı Lab

Hesaplamalı Anlambilim Ders Notları

# Anlamsal Uzayların Metin Sınıflandırmada Kullanımı

Sınıflandırma Algoritması	100 Boyutlu Metinler	10 Boyutlu Metinler
Lineer Regresyonla Sınıflandırma	<b>93.25</b>	89
Pace Regresyonla Sınıflandırma	92.75	88.5
En Yakın Komşu	70	79
Destek Vektör Makineleri	90	89.75
Rastgele Ormanlar (100 ağaçlı)	90.75	87.75

Sınıflandırma Algoritması	Metinlerin Boyut Sayısı	Başarı yüzdesi (tf)	Başarı yüzdesi (tfidf)
Klasik Naive Bayes	4500	87.25	
Diskrit Naive Bayes	4500	85.75	<b>89.25</b>
En Yakın Komşu	4500	34.25	43.5
Destek Vektör Makineleri	4500	87	86.5
Rastgele Ormanlar (100 ağaçlı)	4500	X	X
Lineer Regresyonla Sınıflandırma	4500	X	X
Pace Regresyonla Sınıflandırma	4500	X	X
C4.5	4500	74.75	23.5
Lineer Regresyonla Sınıflandırma	100	81.75	84.5
Pace Regresyonla Sınıflandırma	100	81.5	83.25
En Yakın Komşu	100	71.25	76.25
Destek Vektör Makineleri	100	80.25	87.75
Rastgele Ormanlar (100 ağaçlı)	100	85	84.5



NOVA Research Lab



## Anlamsal Uzayların Metin Sınıflandırmada Kullanımı

- Daha yüksek başarı, Metin temsiliinde daha az hafıza
- Benzerlik matrisi bulurken Google neden kullanılmadı?
- Yöntem dilden bağımsız mı?
- Eşsesli kelimelerin durumu?



### Örnek çalışma-4

Kelime yörüngeleri (Amasyalı, Yener, Kaplan, 2012)

- *İnsanların düşünce süreçleri birbirine ne kadar benzemektedir?*
- *İnsanların düşünce süreçlerini gözlemlemek, direkt ölçmek mümkün değilse de onu, ürünlerinin bazılarıyla (konuşma ve yazılarıyla) dolaylı olarak gözlemlemek mümkündür.*



## Kelime Yörüngeleri

- Yazıları, onları oluşturan kelimelerin zamana göre sıralanmış hali olarak düşünelim.
- Elimizde kelimelerin koordinatları olsa, bu koordinatları yazıdaki sıralarıyla birleştirerek bir yazıyı  $X$  boyutlu bir uzayda bir yörünge olarak ifade edebiliriz.



## Kelime Yörüngeleri

- Aynı kişinin farklı yazılarının yörüngeleri birbirine benzer midir?
- Farklı kişilerin yörüngelerini birbirinden ayırmak mümkün müdür?
- Cevaplar için önce yörüngeleri oluşturmak gerekir. Bunun için ise önce kelime koordinatlarını bulmak gerekir.



## Kelime Yörüngeleri

- Bugün yolda kedi gördüm. Arkasından gittim. Hızlı koşuyordu. Yakalayamadım.
- Hızlı araba kullanmak tehlikelidir. Kediler arabaların altında kalabilirler. Bugün neredeyse eziliyordu bir kedi.



## Yörüngelerin Özellikleri

- 2 özellik:
- Kavramlar arası mesafeler,
- Kavramlar arası açılar.
- $(n, n+1, n+2, n+3)$  arka arkaya gelen 4 koordinat olmak üzere  $(n, n+1)$ ,  $(n, n+2)$ ,  $(n, n+3)$  arası mesafelerin 10'luk histogram değerleri, frekansları ve  $(n, n+1, n+2)$ ,  $(n, n+2, n+3)$  arası açılarının (PI cinsinden) 10'luk histogram değerleri, frekansları



## Yörüngelerin Özellikleri

- $d$  uzunluğundaki bir yörünge de  $d-1$  adet  $(n, n+1)$  arası mesafe ölçülmektedir. Bu  $d-1$  ölçümün eşit aralıklı 10 parçalık histogramı çıkarılmaktadır. Bu histogramın 10 adet değeri ve her değer bir frekansı bulunmaktadır. Dolayısıyla yörüngenin  $(n, n+1)$  arası mesafelerini ifade eden 20 adet özellik çıkarılmaktadır. Bu işlem  $(n, n+2)$ ,  $(n, n+3)$  arası mesafelere ve  $(n, n+1, n+2)$ ,  $(n, n+2, n+3)$  arası açılara da uygulandığında bir metni ifade eden yörüngeye ait 100 özellik bulunmuş olmaktadır



## Deneyisel Sonuçlar

2 yazara ait 35'er yazı,  
Bir yazı kime ait?  
10'lu çapraz geçerleme sonuçları

Algoritma	Sınıflandırma Başarısı (%)
C4.5	97.15
Naive Bayes	97.15
En yakın komşu	95.71
Destek Vektör Makineleri	98.57

2 yazara ait 70 yazının  $(n, n+1)$  arası mesafelerinin histogramının ilk parçasının frekanslarının histogramı

Her yazara ait 35'er yazı,  
Bir yazı kime ait?  
10'lu çapraz geçerleme sonuçları

Algoritma	3 yazar	4 yazar	8 yazar
C4.5	91.43	90.71	53.21
Naive Bayes	95.24	92.14	62.5
En yakın komşu	80.95	62.14	50
Destek Vektör Makineleri	94.29	85.71	68.57
Zero Rule	33.33	25	12.5



## Sonuç ve Tartışma

- 2 kişiye ait 35'er yazıdan oluşan veri kümesinde bir metnin yazarını tanıma başarısı %98
- Yazar sayısı arttıkça başarının düşüyor (8 yazar için %68)
- Metinde kullanılan kavramlardan bağımsız olması (sadece yörünge'nin özelliklerinin kullanılıyor olması) aynı kişinin farklı konularda yazdıklarından elde edilen yörüngeler benzer



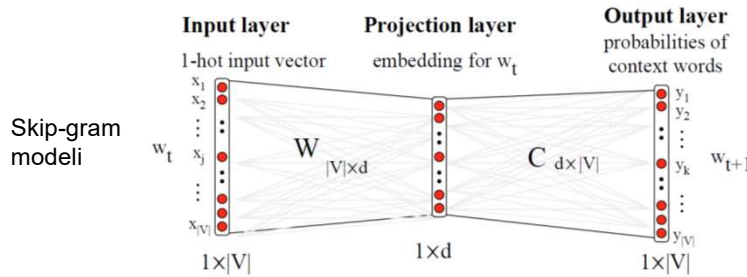
## Olası kullanım alanları

- Psikolojik hastalıkların tespiti
- Psikolojik hastalıkların düşünce süreçleri üzerindeki etkilerinin araştırılması
- Cinsiyet, yaş, eğitim farklılıklarının düşünce süreçleri üzerindeki etkilerinin araştırılması



## Sinirsel Dil Modelleri

- Word2vec\*, fasttext\*\*
- Tek gizli katmana sahip bir yapay sinir ağı
- $C=2$  için bağlam =  $(w_{t-2}, w_{t-1}, w_t, w_{t+1}, w_{t+2})$
- Skip-gram: giriş:  $w_t$ , çıkış:  $(w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2})$  in herbiri
- C-bow: giriş:  $(w_{t-2}, w_{t-1}, w_{t+1}, w_{t+2})$  in herbiri, çıkış:  $w_t$



\* 2013, T. Mikolov ve ark. "Distributed representations of words and phrases and their compositionality"

\*\* 2017, P. Bojanowski ve ark. "Enriching word vectors with subword information"

NOVA Research Lab

## Sinirsel Dil Modelleri

- $W$  matrisi kelime vektörleri
- Başlangıçta rasgele atanırlar. Eğitim bittiğinde kullanıma hazırlar.
- Amaç: bir kelimenin ( $w$ ) genelde etrafında bulunan kelimelerle ( $c$ ) benzer vektörlere, diğerleriyle ( $w_i$ ) farklı vektörlere sahip olması
- 2 vektör benzer ise iç çarpımları yüksektir.
- Negatif örnekleme Skip-gram için hata fonksiyonu:

$$\log \sigma(c \cdot w) + \sum_{i=1}^K \mathbb{E}_{w_i \sim p(w)} [\log \sigma(-w_i \cdot w)]$$

$$\sigma(x) = \frac{1}{1+e^x}$$



<https://medium.com/@Aj.Cheng/word2vec-3b2cc79d674>

<https://web.stanford.edu/~jurafsky/slp3/slides/vector2.pdf>

NOVA Research Lab

## Eğitim kümesi oluşturmak

- Text:  $w_1 w_2 w_3 w_4 w_5 w_6 w_7 \dots$
- Skip gram için eğitim örnekleri (giriş, çıkış)  $C=2$  için:
  - benzer olmasını istediğimiz ikililer:  $(w_3, w_1), (w_3, w_2), (w_3, w_4), (w_3, w_5)$
  - farklı olmasını istediğimiz ikililer:  $(w_3, t_1), (w_3, t_2), (w_3, t_3) \dots (w_3, t_K)$  burada  $t_1 \dots t_K$   $p(w_i)$  ye göre rasgele seçilen kelimeler
- Sonra  $w_4$  ile devam ...
- C-bow için ikilileri tersine çevir



## Fasttext

- Birlikte bulunmaya ek olarak kelime içi n-gram benzerliklerini (sub-word similarity) de dikkate alır.
- Kelimeleri n-gram larına parçalar ve her bir n-gram içinde vektörler bulur.
- Bu sayede sözlükte olmayan (test örneklerinde çok karşılaşılan bir durum) kelimelerin vektörlerine bulur.
- Türkçe gibi ek alan diller için avantajlı (aynı kelimenin farklı ekler almış halleri, yanlış yazılmış halleri birbirine yakın koordinatlara sahip olur)
- Ama daha yavaş



## kodlar

- word2vec\_modeli\_hazir.m (kullanım)
- word2vec\_modeli\_egitim.m (eğitim)
- emb42\_ngramsiz.mat (kelime vektörleri)
- fast\_emb.mat (kelime vektörleri)
- haber42bin.mat (derlem)



## Hazır fasttext\* vektörleri ile ilk satırdaki kelimelere en yakın 19 kelime

"cimbom"	"terlik"	"virüs"	"elma"	"aslan"	"koş"	"bırak"
"cimbombom"	"terliksi"	"virüsün"	"elma,armut"	"aslaner"	"koşo"	"bıraktın"
"lik--"	"terlikler"	"virüslü"	"armut"	"aslana"	"koşoy"	"bırakın"
"penaltılarla"	"terlikleri"	"virüsle"	",armut"	"aslani"	"koşaç"	"bırakmam"
"#beşiktaş"	"noterlik"	"virüse"	"şeftali"	"aslankeser"	"koşa"	"bırakmam"
"rövanşta"	"kasiyerlik"	"virüstür"	"meyveş"	"aslani"	"koşagaç"	"bırak"
"lik/"	"önlük"	"virüsü"	"meyva"	"aslanin"	"koşovi"	"bıraksın"
"penaltısonuç"	"yeterlik"	"virüsler"	"elmaağacı"	"aslanyürek"	"koşun"	"bıraktım"
"beşiktaşlı"	"dülgerlik"	"virüsün"	"kayısı"	"aslanbaş"	"koşamaz"	"bırakayım"
"cimbilli"	"dizlik"	"virüsten"	"meyvelik"	"aslanhan"	"koşalı"	"bırakırsan"
"gülesin"	"göynek"	"virüste"	"ayva"	"özaslan"	"koşacı"	"bıra"
"galatasaray"	"yerlik"	"virüsleri"	"üzüm"	"aslanbay"	"koşaca"	"bırakırsam"
"üsk"	"ayakkabıcılık"	"virüslere"	"elmakaya"	"aslanduz"	"koşköpür"	"bırakmayın"
"eledim"	"tonmaysterlik"	"mimivirüs"	"buğday,şeker"	"aslanpençesi"	"koşay"	"bıraktık"
"fenerbahçe"	"poşu"	"virütik"	"badem"	"aslanıtuğ"	"koşan"	"bırakmalı"
"haftabjk"	"köynek"	"virüslerde"	"arpa,buğday"	"aslanla"	"koşabat"	"bırakalım"
"beşiktaş"	"yakıcılık"	"virüsleri#grup"	"armutu"	"aslanları"	"beliy"	"bırakalı"
"rövanşında"	"rençberlik"	"poliovirüs"	"şeftalinin"	"sanaslan"	"koşanlar"	"bırakmadım"
"cim"	"terli"	"virüslerin"	"elmax"	"aslanlar"	"koşunca"	"bıraksam"
"kupađaki"	"üzerlik"	"virüslerle"	"armudu"	"yavuzaslan"	"koşkin"	"bırakma"



\*Fasttext vektörleri 416 bin kelime içermektedir.  
Ngram benzerliği de kullanır.



## 42 bin haber metniyle eğitilen vektörler\*, ngram benzerliği yok ilk satırdaki kelimelere en yakın 19 kelime

"cimbom"	"istanbul"	"virüs"	"elma"	"aslan"	"ölmek"	"bırak"
"kırmızılılar"	"sergi"	"virüsü"	"meyve"	"aslanın"	"yaşamak"	"bırakın"
"beşikta"	"gerçekleşecek"	"virüsün"	"meyveler"	"lakaplı"	"istemiyorum"	"dedikleri"
"wesley"	"sahipliğinde"	"bulaşan"	"lif"	"cimbom"	"asker"	"kardeşim"
"didier"	"katılımıyla"	"hiv"	"sebze"	"kırmızılı"	"erbakan"	"üzü"
"galatasaray"	"fashion"	"hastalık"	"limon"	"kırmızılılar"	"chavez"	"git"
"transfere"	"konser"	"yoluyla"	"çürük"	"serhat"	"lütfe"	"dediğini"
"kırmızılı"	"buluşacak"	"grip"	"salata"	"devler"	"bekleyen"	"hayattan"
"kırmızılıların"	"sarayında"	"aşı"	"pirinç"	"didier"	"kanserin"	"yahu"
"eboue"	"inci"	"zararlı"	"buğday"	"muhammet"	"ameliyata"	"genişletilmiş"
"cim"	"sergisi"	"yazılımı"	"zeytinyağı"	"toure"	"anın"	"memnunum"
"sarıkırmızılı"	"cnr"	"bilgisayar"	"patates"	"drogbay"	"rio"	"aşkına"
"devler"	"seçkin"	"solunum"	"çiğ"	"lokomotiv"	"aileme"	"temsilcimiz"
"aysal"	"global"	"enfeksiyonları"	"muz"	"payını"	"ölümü"	"etkisine"
"ezeli"	"modern"	"kullanıcının"	"ürünlerde"	"ulaş"	"istemediğini"	"özette"
"terimin"	"festival"	"belirtileri"	"findık"	"oya"	"ağabeyim"	"dedim"
"ligindeki"	"louis"	"ağrısına"	"sarımsak"	"bulut"	"koydular"	"açıklamalarda"
"deplasmanında"	"dünyaca"	"enfeksiyon"	"portakal"	"yitirdi"	"özbeğenin"	"gidin"
"futbolcunun"	"konseri"	"bulaşıcı"	"hapı"	"ceren"	"yasak"	"yazdı"
"beraberlik"	"tarihleri"	"saldırganlar"	"ihracatı"	"ekrem"	"söylemiş"	"oynuyor"



\*Vektörler 19 bin kelime içermektedir.

NOVA Research Lab

## Kelime vektörlerinin başarısını nasıl ölçeriz?

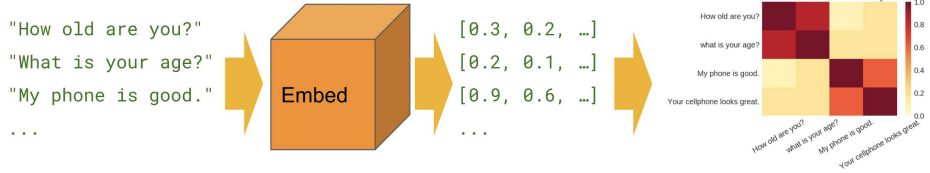
- Dahili / direkt yöntemler
  - İnsanların verdikleri cevaplarla korelasyon
  - Analoji veri kümeleri:
    - [https://aclweb.org/aclwiki/Google\\_analogy\\_test\\_set\\_\(State\\_of\\_the\\_art\)](https://aclweb.org/aclwiki/Google_analogy_test_set_(State_of_the_art))
    - <https://www.aclweb.org/anthology/C16-1332.pdf>
- Harici / dolaylı yöntemler
  - Çeşitli NLP uygulamalarında kullanımları üzerinden



NOVA Research Lab

## Cümle ve metin vektörleri

- En basit yol: içindeki kelimelerin ortalaması
- Universal sentence encoder  
<https://arxiv.org/abs/1803.11175>



## Kelime Öbekleri

- 2, 3, 4 uzunluğunda genelde birlikte kullanılan kelime birlikleri içerdikleri her bir kelimedenden ayrı bir kavramı ifade edebilirler.
- “Fatih Sultan Mehmet” ifadesi, “Fatih”, “Sultan” ve “Mehmet” in her birinden, farklı bir kavramı ifade eder.



## Kelime Öbekleri

- Bu ifadeleri nasıl buluruz?
  - En basiti: frekansı en yüksek olanlar
  - Kelime öbeği olma olasılığı (ti):
    - ti nin frekansı/ti nin içerdiği tekil kelimelerin frekanslarının minimumu
    - $\log_2(\text{ti nin frekansı} / \text{ti nin içerdiği tekil kelimelerin frekanslarının çarpımı})$
- En basitini deneyelim.



## Kod kelime\_obekleri.m

- 42bin haber metni içindeki 1, 2, ve 3 uzunluğunda kelime ngramlarını bulalım
- Nadir kelimeleri silmezsek 11 milyon ngram çıkıyor ☹️
- 50 den az geçen kelimeleri silersek
  - kelime silme öncesi tekil kelime sayısı= 356.680
  - kelime silme öncesi dokümanlardaki toplam kelime sayısı= 9.986.668
  - silinen tekil kelime sayısı= 336.019
  - kelime silme sonrası dokümanlardaki toplam kelime sayısı= 8.320.579

```
Counts: [41992x8141753 double]
Vocabulary: [1x20661 string]
Ngrams: [8141753x3 string]
NgramLengths: [1 2 3]
NumNgrams: 8141753
NumDocuments: 41992
```

→ Ngram sayısı 8 milyon



## Hesaplama Anlambilim Ders Notları

### Frekansı en yüksek 1 2 3 kelime ngramlar

Ngram			Count	NgramLength					
					Ngram			Count	NgramLength
"ve"	"	"	1.8454e+05	1					
"bir"	"	"	1.4992e+05	1					
"ı"	"	"	1.2161e+05	1					
"bu"	"	"	93130	1	"türkiye"	"ı"	"nin"	2187	3
"da"	"	"	73033	1	"türkiye"	"ı"	"de"	2171	3
"de"	"	"	70376	1	"ı"	"diye"	"konuştu"	2168	3
"için"	"	"	57285	1	"spor"	"toto"	"super"	1580	3
"ile"	"	"	54426	1	"şöyle"	"devam"	"etti"	1509	3
"çok"	"	"	39307	1	"başbakan"	"recep"	"tayyip"	1407	3
"olarak"	"	"	34066	1	"türkiye"	"ı"	"ye"	1017	3
					"recep"	"tayyip"	"erdoğan"	996	3
					"yönetim"	"kurulu"	"başkanı"	983	3
					"bir"	"kez"	"daha"	909	3
					"şunları"	"soyledi"	"ı"	744	3
					"i'stambul"	"ı"	"da"	728	3
					"başt"	"olmak"	"üzere"	720	3
					"toto"	"super"	"ligde"	717	3
					"genel"	"başkan"	"yardımcısı"	708	3
					"çok"	"önemli"	"bir"	633	3
"ı"	"nin"	"	9411	2	"erdoğan"	"ı"	"ın"	627	3
"ı"	"ın"	"	9061	2	"şöyle"	"konuştu"	"ı"	576	3
"ı"	"de"	"	8796	2	"olduğunu"	"ifade"	"eden"	561	3
"ı"	"da"	"	8793	2	"aa"	"muhabirine"	"yaptığı"	559	3
"ı"	"in"	"	7087	2	"bir"	"süre"	"sonra"	556	3
"ı"	"ın"	"	6578	2	"muhabirine"	"yaptığı"	"açıklamada"	555	3
"türkiye"	"ı"	"	6541	2	"bir"	"an"	"once"	551	3
"diye"	"konuştu"	"	6450	2	"olay"	"yerine"	"gelen"	526	3
"ı"	"dedi"	"	6240	2					
"ya"	"da"	"	5226	2					

NOVA Research Lab

## Hesaplama Anlambilim Ders Notları

- 8 milyon ngramdan sadece 36 bin'i 50 den fazla kez geçiyor. Sadece bunlarla ilerleyelim. Tekil kelime sayısı 20 bin.
- Ngramları kendi içlerindeki kelimelerden ayrı kendi başlarına anlamları / koordinatları olsun diye tek kelimelere çevirelim.

"bitirdikten sonra"	"bitirdikten sonra"
"işlemlerinde"	"işlemlerinde"
"ısrarlı"	"ısrarlı"
"dağıldı"	"dağıldı"
"fırtınası"	"fırtınası"
"yer alırken"	"yer alırken"
"700 milyon"	"700 milyon"
"sırada yer alırken"	"sıradayeralırken"
"telefonundan"	"telefonundan"
"yi'ne"	"yi'ne"
"kız kardeşi"	"kızkardeşi"
"cep telefonundan"	"ceptelefonundan"
"i'nternet üzerinden"	"i'nternetüzerinden"
"başbakan da"	"başbakanda"
"ve şöyle"	"veşöyle"
"bir tane"	"birtane"

NOVA Research Lab

## Haberlerin yeni hali

- (1,1) 111 tokens: ortak vizyonumuz var **dışişleribakanı** davutoğlu yunanistan **iletürkiye** arasındaki farklılıkların ortak vizyon ile çözülebileceğini söyledi **dışişleribakanıahmet davutoğlu**ve yunan mevkidaşı dimitris avramopulos başbaşa ve heyetler arası **görüşmelerinardından ortakbasıntoplantısı** düzenledi davutoğlu yunanistanda kendisine gösterilen sıcak karşılama **içinteşekkür** ederek ziyaretinin **türkyunan** dostluğunun bir nişanesi **olduğunudile** getirdi **ikiülke** arasında görüş ayrılıkları **vefarklı** yaklaşımlar olabileceğini **amaaynı** vizyonun paylaşılmasıyla aradaki sorunların problemlerin **daharahat** çözülebileceğini **ifadeeden** bakan davutoğlu samaras ve avramopulos ile görüşmelerinden **eldeettiği** sonucun **dabu** yönde olduğunu söyledi davutoğlu **ikiyılönce** ydiki bu anlayışla kurduklarını yunanistanda **istikrarlıbir** hükümet kurulur kurulmaz da toplantının ikincisini gerçekleştirmeyi planladıklarını **vebugün** ön hazırlıkları **önemliölçüde** tamamladıklarını söyleyerek birçok alanda ortak mutabakatların ilk tohumları atıldı ...



42 bin haber metniyle eğitilen kelime ngramı vektörler\*, karakter ngram benzerliği yok ilk satırdaki kelimelere en yakın 19 kelime

"cimbom"	"istanbul"	"virüs"	"elma"	"aslan"	"ölmek"	"bırak"
"galatasaray"	"sergi"	"virüsün"	"sebze"	"aslanın"	"giymek"	"git"
"g"	"gerçekleşecek"	"virüsü"	"soğan"	"hakan"	"direniş"	"dedim"
"saray"	"festival"	"bulaşan"	"meyve"	"pelin"	"benim"	"başkabilirsey"
"galatasarayın"	"fashion"	"bilgisayara"	"domates"	"irem"	"okumak"	"gel"
"sarıkırmızılılar"	"istanbul"	"yazılımı"	"patates"	"sevgilisiyle"	"cenazetöreni"	"bul"
"drogba"	"filarmoni"	"bulaşıcı"	"çürük"	"yalçınkaya"	"üzereyken"	"yahu"
"terim"	"konser"	"hastalık"	"salata"	"mehmet"	"benim"	"demesiüzereine"
"fenerbahçe"	"fuari"	"bilgisayar"	"portakal"	"kaplı"	"ölüm"	"hiç"
"sarıkırmızılı"	"lütifikırdar"	"bakteri"	"tavuk"	"yaşar"	"şilri"	"diyecek"
"realmadrid"	"festivali"	"enfeksiyon"	"meyve"	"naz"	"yaşamak"	"kızım"
"snejder"	"solo"	"bağışıklık"	"tadı"	"neslihan"	"öldürmekle"	"geridönen"
"i"	"orkestrasi"	"kullanıcıların"	"kabak"	"yusu"	"hastalığınedeniyle"	"mahsun"
"fildişli"	"sanatseverlerle"	"antibiyotik"	"ürünlerin"	"abdurrahman"	"istemiyorum"	"görmüyorum"
"sarıkırmızılıların"	"konseri"	"şebeke"	"muz"	"bom"	"sormak"	"ama"
"galatasarayda"	"bilkent"	"tedavisinde"	"üzüm"	"hasan"	"geldin"	"müddetçe"
"beşiktaş"	"grand"	"telefonlar"	"şekerli"	"erol"	"bana"	"bilmiyorum"
"ligi"	"prixsi"	"dolandırıcılık"	"sebze"	"2523"	"ogün"	"sana"
"cim"	"caz"	"grip"	"kiraz"	"çini"	"çocuğumu"	"sapan"
"şampiyonlar"	"konserler"	"porno"	"meyveler"	"shanghai"	"evlenmek"	"hiçkimse"



\*Vektörler 26 bin kelime/ngram içermektedir.

## Faydalı bağlantılar

- Dil modellerini değerlendirmek: <https://thegradient.pub/understanding-evaluation-metrics-for-language-models/>
- Klasik dil modelleme: <https://github.com/vrann/next-word-prediction>
- Generalization through Memorization: Nearest Neighbor Language Models: <https://arxiv.org/abs/1911.00172>

