



BTK
AKADEMİ

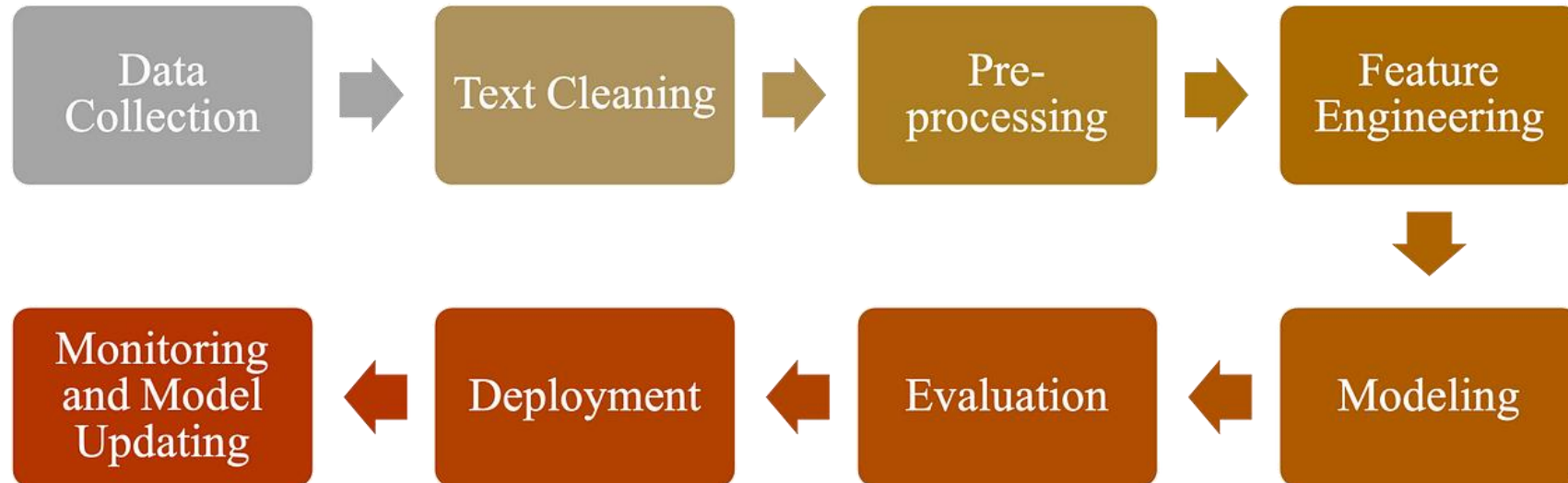
Bölüm 5

Metin Temsili

Metin Temsili

- Metin temsili, bir metni sayısal veya başka türde bir formatta temsil etme işlemidir.

NLP Pipeline



https://alvinntnu.github.io/NTNU_ENC2045_LECTURES/nlp/nlp-pipeline.html

Metin Temsili Neden Yapılmalıdır?

- Bilgisayarların Anlayabilmesi
- Öznetelik Çıkartma
- Model Eğitim

Metin Temsili Yöntemleri

- Bag of Words (BoW)
- TF-IDF (Term Frequency-Inverse Document Frequency)
- N-Gram Modelleri
- Word Embeddings
- Transformers Tabanlı Metin Temsili

Bag of Words (BoW)

Bag of Words Nedir?

- Bag of Words (BoW), doğal dil işleme (NLP) ve metin madenciliğinde kullanılan temel bir metin temsili yöntemidir.
- BoW, metinlerdeki kelimeleri sayısal verilere dönüştürür ve metinlerin analizini sağlar.



<https://dataaspirant.com/bag-of-words-bow/>

BoW Yönteminin İşleyişi

- Kelime Kümesi Oluşturma
- Kelime Frekansı Hesaplama
- Vektör Temsili

BoW Yönteminin İşleyişi

- **Kelime Kümesi Oluşturma**

Metin 1: Kedi evde

Kelime Kümesi: ["Kedi", "evde", "bahçede"]

Metin 2: Kedi bahçede

- **Kelime Frekansı Hesaplama**

Metin 1: "Kedi evde"

"Kedi": 1

"evde": 1

"bahçede": 0

Metin 2: "Kedi bahçede"

"Kedi": 1

"evde": 0

"bahçede": 1

- **Vektör Temsili**

Metin 1: [1, 1, 0]

Metin 2: [1, 0, 1]

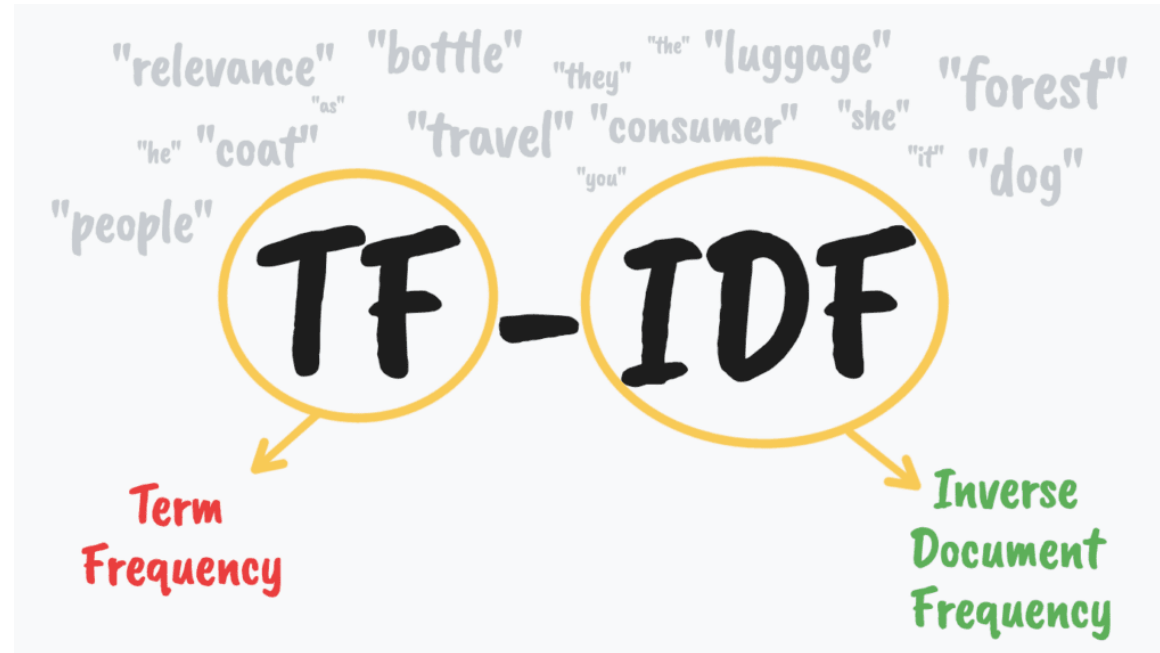


BTK
AKADEMİ

TF-IDF (Term Frequency-Inverse
Document Frequency)

TF-IDF (Term Frequency-Inverse Document Frequency)

- TF-IDF (Term Frequency-Inverse Document Frequency), metin madenciliğinde ve bilgi erişiminde sıkça kullanılan bir özellik çıkarım yöntemidir.
- TF-IDF, kelimelerin belgeler içinde ne kadar önemli olduğunu belirlemek için kullanılır.



TF-IDF (Term Frequency-Inverse Document Frequency)

- **Term Frequency (TF):** Bir kelimenin bir belgede ne kadar sık geçtiğini ölçer.
- **Inverse Document Frequency (IDF):** Bir kelimenin tüm belgelerdeki yaygınlığını ölçer. Bir kelimenin çok belgede geçiyorsa, o kelime çok fazla bilgi sağlamaz.

$$TF(t, d) = \frac{\text{Sayac}(t, d)}{\text{Toplam Kelime Sayısı}(d)}$$

Burada, $\text{Sayac}(t, d)$ kelimenin t bir belgede d sayısıdır.

$$IDF(t, D) = \log \left(\frac{\text{Toplam Belgeler Sayısı}(D)}{1 + \text{Belgelerde Geçen Sayısı}(t)} \right)$$

Burada, $\text{Belgelerde Geçen Sayısı}(t)$ kelimenin t belgelerdeki sayısıdır.

$$TF-IDF(t, d, D) = TF(t, d) \times IDF(t, D)$$



BTK
AKADEMİ

N-Gram Modelleri

N-Gram Modelleri Nedir?

- Bir dil modelinde kullanılan kelime veya karakter dizisinin uzunluğunu belirten bir terimdir.
- N-Gram modelleri, metinleri n kelimelik veya n karakterlik kısımlara bölerek analiz eder.

This is Big Data AI Book

Uni-Gram

This	Is	Big	Data	AI	Book
------	----	-----	------	----	------

Bi-Gram

This is	Is Big	Big Data	Data AI	AI Book
---------	--------	----------	---------	---------

Tri-Gram

This is Big	Is Big Data	Big Data AI	Data AI Book
-------------	-------------	-------------	--------------

N-Gram Modelleri

- "Bu bir örnek metindir"
- **Unigram (n=1)**
 - ['Bu', 'bir', 'örnek', 'metindir']
- **Bigram (n=2)**
 - ['Bu bir', 'bir örnek', 'örnek metindir']
- **Trigram (n=3)**
 - ['Bu bir örnek', 'bir örnek metindir']

N-Gram Modelleri Kullanım Alanları

- Metin Modelleme
- Metin Sınıflandırma
- Metin Üretimi
- Metin Benzerliği

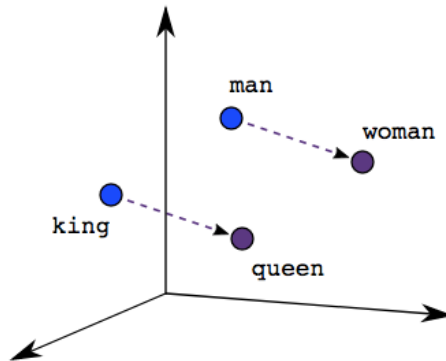


BTK
AKADEMİ

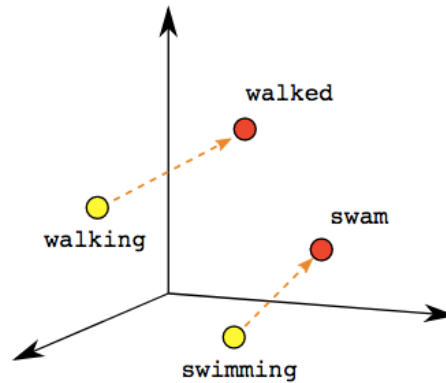
Word Embeddings

Word Embeddings

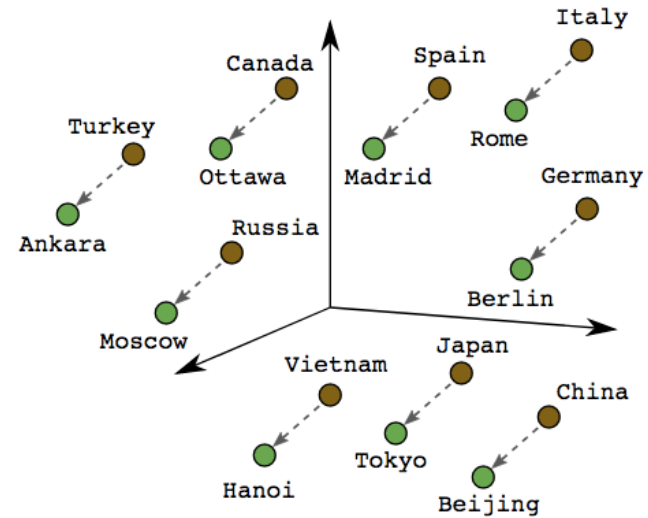
- Word Embeddings (Kelime Gömmme), doğal dil işleme (NLP) ve makine öğreniminde kullanılan bir tekniktir.
- Kelimeleri, genellikle sürekli bir vektör uzayında anlamlı temsil edecek şekilde sayısal vektörlere dönüştürür.
- Bu temsiller, kelimeler arasındaki anlamsal ve dilbilgisel ilişkileri yakalamayı hedefler.



Male-Female



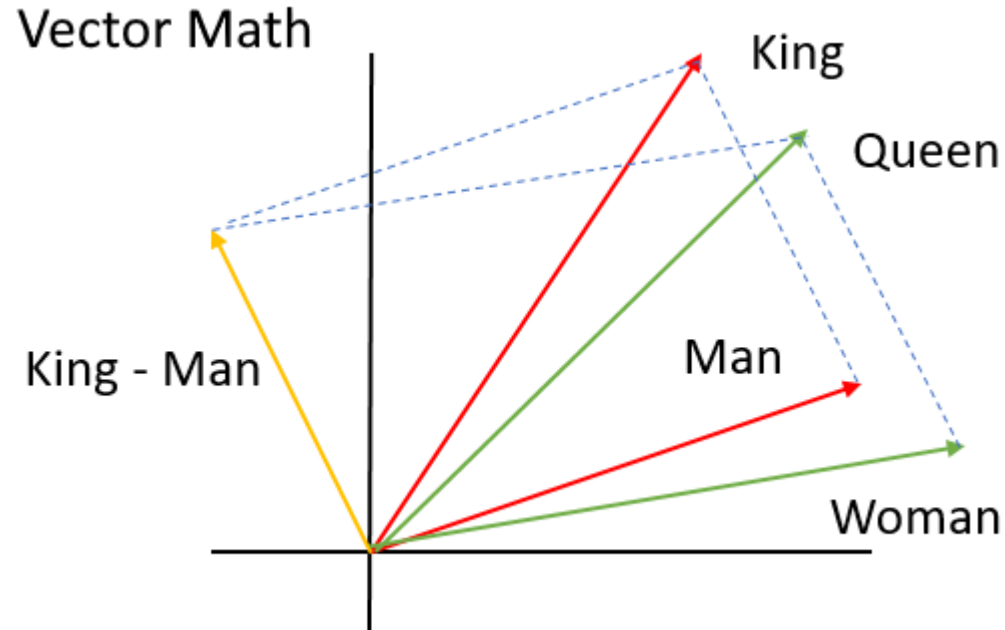
Verb Tense



Country-Capital

Word Embeddings Özellikleri

- Anlamsal Benzerlik
 - Örneğin, "king" ve "queen" kelimeleri benzer vektörler alabilir.
- Matematiksel İşlemler
 - Örneğin, "king" - "man" + "woman" = "queen" hesaplaması yapılabilir.
- Kapsamlılık



Word Embeddings Modelleri

- **Word2Vec:** Google tarafından geliştirilen, kelimeleri vektörlere dönüştüren ve bu vektörleri dildeki ilişkileri yakalayacak şekilde eğiten bir modeldir.
- **GloVe (Global Vectors for Word Representation):** Stanford Üniversitesi tarafından geliştirilen, kelime gömme temsillerini kelime ortaklıklarını yakalayacak şekilde hesaplayan bir modeldir.
- **FastText:** Facebook tarafından geliştirilen ve kelime gömme temsillerini kelime alt-birimlerini de dikkate alarak hesaplayan bir modeldir.



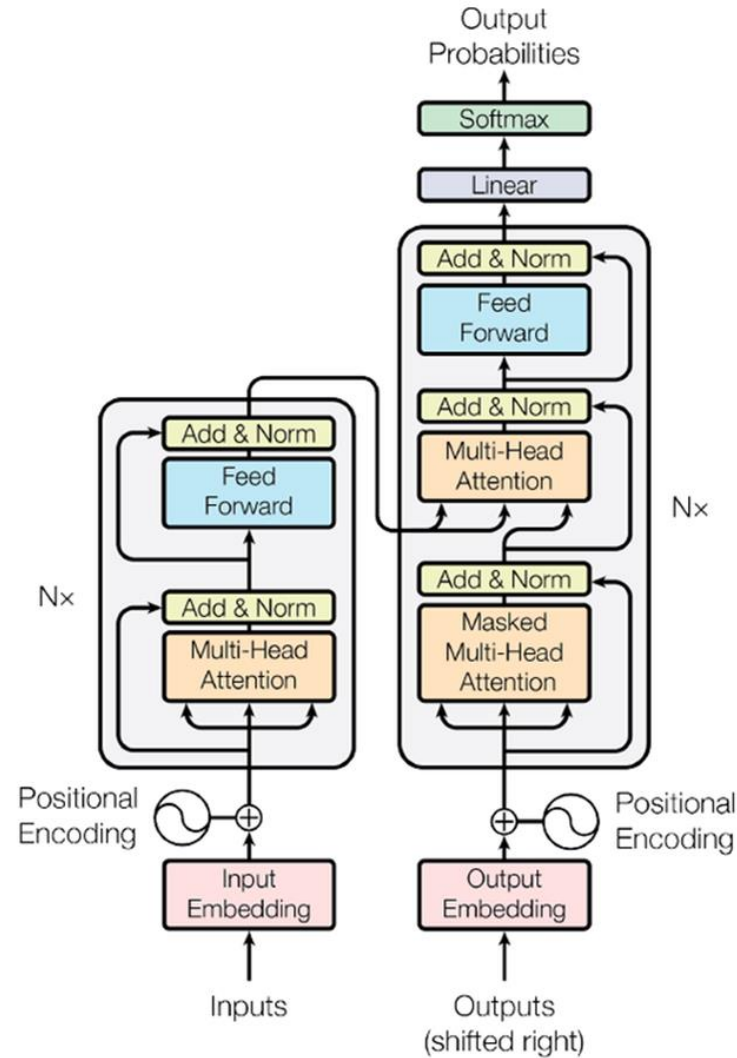
BTK
AKADEMİ

Transformers Tabanlı Metin Temsili

Transformers Tabanlı Metin Temsili

- Transformers, doğal dil işleme (NLP) ve diğer yapay zeka alanlarında son yıllarda devrim niteliğinde yenilikler getiren bir mimaridir.
- İlk olarak 2017 yılında Google tarafından yayınlanan "Attention is All You Need" adlı makalede tanıtılmıştır.
- Neden Transformers
 - Bağlamı Daha İyi Anlama
 - Paralel İşleme Yeteneği
 - Çeşitli NLP Görevlerinde Kullanım
 - Önceden Eğitilmiş Modellerin Yeniden Kullanımı
- En Bilindik Transformers Modelleri
 - BERT (Bidirectional Encoder Representations from Transformers)
 - GPT (Generative Pre-trained Transformer)

Transformers Tabanlı Metin Temsili



Transformers Tabanlı Metin Temsili

- **Attention**, modelin belirli girdi parçalarına farklı derecelerde dikkat göstermesine olanak tanır.
- Özellikle, bir kelimenin diğer kelimelerle olan ilişkisini anlamak için kullanılır.
- Örneğin, "Kedi hızlıdır" cümlesinde, "hızlıdır" kelimesi "Kedi" kelimesine olan dikkat skorlarını hesaplar:
 - **Sorgu (Query) ve Anahtar (Key) Çarpımı:** "Kedi" kelimesinin "hızlıdır" kelimesi ile olan ilişki skoru hesaplanır.
 - **Dikkat Skoru:** Bu skora göre "Kedi" kelimesinin temsili günceller.

Transformers Tabanlı Metin Temsili

- **Input Embedding**, girdi verilerini modelin işleyebileceği bir formata dönüştürmek için kullanılan bir tekniktir.
- Örnek: Bir cümle düşünelim: "Kedi hızlıdır.«
- Kelime Vektörleri: Bu cümledeki kelimeler, word2vec, GloVe veya BERT gibi bir embedding tekniği kullanılarak sayısal vektörlere dönüştürülür.
 - "Kedi" \rightarrow [0.21, -0.32, 0.87, ...]
 - "hızlıdır" \rightarrow [-0.13, 0.45, -0.20, ...]

Transformers Tabanlı Metin Temsili

- **Multi-Head Attention**, attention mekanizmasının birden fazla başlıkla (head) çalıştığı bir tekniktir.
- Bir cümledeki her kelime, diğer kelimelerle olan ilişkilerini farklı açılardan öğrenmek isteyebilir.
- Örneğin, "Kedi" kelimesinin "hızlıdır" kelimesiyle ilişkisini anlamak için birden fazla dikkat başlığı kullanılır.
 - Başlık 1: "Kedi" ve "hızlıdır" arasındaki anlam ilişkisini öğrenir.
 - Başlık 2: "Kedi" ve "hızlıdır" arasındaki gramatik ilişkileri öğrenir.
 - Başlık 3: "Kedi" kelimesinin cümledeki konumunu öğrenir.

Transformers Tabanlı Metin Temsili

- **Masked Multi-Head Attention**, modelin gelecekteki kelimeleri görmesini engeller, yani model sadece geçmiş bilgileri kullanarak tahminde bulunur.
- Örneğin, "Kedi _____ hızlıdır" cümlesinde, model "Kedi" ve "hızlıdır" arasındaki ilişkilere dayanarak boşluğa "oldukça" kelimesini tahmin eder.

Transformers Tabanlı Metin Temsili

- **Add & Norm**, bir katman çıktı ile giriş arasındaki kısa yolu (residual connection) ekleyip ardından layer normalization uygulayan bir adımdır.
- **Feed-Forward Network**, her encoder ve decoder katmanında bulunan bir ağıdır.
- **Output Embedding**, modelin çıktısını temsil eden ve genellikle bir dil modelinde kullanılan bir tekniktir.

Metin Temsili Yöntemlerinin Karşılaştırılması

Yöntem	Temel Özellikler	Kullanım Kolaylığı	Sonuçların Başarı Durumu
Bag of Words (BoW)	Kelime frekanslarına dayalı, sıklık matrisleri oluşturur.	Basit, doğrudan uygulanabilir.	Genellikle düşük, bağlam bilgisinden yoksundur.
TF-IDF	Kelime sıklığına ek olarak, kelimenin belgelerdeki önemini ölçer.	Kolay, standart kütüphaneler mevcut.	Orta, bağlam bilgisi kısıtlıdır ama bilgiye değer katar.
N-grams	Kelime ya da karakter n-gramlarını kullanarak bağlamı yakalar.	Orta, işlem gücü ve bellek kullanımı artar.	Orta, bağlam bilgisi artırılabilir ancak model karmaşıklığı artar.
Word Embeddings (GloVe, Word2Vec, FastText)	Kelimeleri vektörlere dönüştürür, anlam ilişkilerini yakalar.	Orta, önceden eğitilmiş modeller mevcut.	Yüksek, bağlamı daha iyi anlar, semantik ilişkiler sağlar.
Transformers (BERT, GPT-3, vb.)	Derin öğrenme temelli, bağlamı dikkat mekanizması ile yakalar.	Orta-İleri, genellikle yüksek hesaplama gücü gerektirir.	Çok yüksek, bağlamı derinlemesine anlar, çeşitli NLP görevlerinde başarılı.