

Sağlık Alanında Büyük Dil Modellerinin Adaptasyonu ve Veri Toplama ve Hazırlama Süreci: Bir Pratik Örnek

Özet:

Büyük dil modelleri (BDM), tıbbi bilgiye erişimi iyileştirmek, hastalarla iletişimi güçlendirmek ve yeni tedaviler geliştirmek gibi potansiyelleri ile sağlık alanında devrim yaratma potansiyeline sahiptir. Ancak, BDM'lerin sağlık alanında etkili ve güvenli bir şekilde kullanılabilmesi için yüksek kaliteli veri toplamak ve hazırlamak gereklidir. Bu makalede, doktorlar tarafından hastaların sorduğu sorulara cevaplar verilen ve herkese açık olarak paylaşılan bir web sitesinden elde edilen doktor anonim profilleri ve soru-cevap verileri kullanılarak BDM adaptasyonu için veri toplama ve hazırlama süreci ele alınmıştır. Veri toplama, veri birleştirme, boş değerlerin işlenmesi, veri tipi dönüşümü, veri temizleme, veri kalite kontrolü, BDM eğitime hazırlık ve metin ön işleme gibi adımlar detaylı olarak açıklanmıştır. Hazırlanan veriler kullanılarak Meta şirketinin geliştirdiği LLAMA 3 modeli ve YTÜ COSMOS yapay zeka araştırma grubunun geliştirdiği cosmosGPT v0.1 modeli üzerinde fine-tuning işlemi gerçekleştirilmiştir. Böylece modeller, sağlık alanında Türkçe sorulan sorulara daha iyi cevaplar verebilmeye başlamıştır. Bu çalışma, BDM'lerin sağlık alanında kullanımı için kaliteli veri toplama ve hazırlamanın önemini vurgulamaktadır.

Anahtar Kelimeler: Büyük Dil Modelleri, Sağlık Alanı, Veri Toplama, Veri Hazırlama, Kalite Kontrol, BDM Eğitimi, Fine-tuning, LLAMA 3, cosmosGPT.

1. Giriş:

Sağlık hizmetleri, insan yaşamının en önemli ve hassas alanlarından biridir. Sağlıklı bir yaşam sürdürebilmek için doğru ve zamanında tıbbi bilgiye erişim, hastalıkların doğru teşhisi ve etkili tedavi yöntemlerinin uygulanması büyük önem taşımaktadır. Son yıllarda teknolojinin hızla gelişmesi, sağlık hizmetlerinin sunumunda, teşhis ve tedavi süreçlerinde yeni ve heyecan verici olanaklar yaratmaktadır. Yapay zeka (YZ), bu dönüşümün ön saflarında yer alan teknolojilerden biridir. YZ, karmaşık tıbbi verileri analiz ederek, hastalıkları teşhis etmede, kişiselleştirilmiş tedavi planları oluşturmada ve hatta yeni ilaçlar keşfetmede kullanılabilme potansiyeline sahiptir.

Bu potansiyelin en önemli temsilcilerinden biri de Büyük Dil Modelleri'dir (BDM). BDM'ler, devasa metin veri kümeleri üzerinde eğitilmiş derin öğrenme algoritmalarıdır. Bu modeller, doğal dili anlama, yorumlama ve üretme konusunda son derece yeteneklidirler. İnsanlar gibi metinleri okuyabilir, yazabilir, özetleyebilir ve hatta farklı diller arasında çeviri yapabilirler.

Sağlık alanında BDM'lerin potansiyel uygulamaları oldukça geniştir. Örneğin:

- **Tıbbi Bilgiye Erişim:** Hastalar, BDM'ler aracılığıyla tıbbi bilgilerine kolayca erişebilir, hastalıkları hakkında bilgi edinebilir, semptomlarını değerlendirebilir ve tedavi seçenekleri hakkında bilgi alabilirler.
- **Hasta-Doktor İletişimi:** BDM'ler, hasta-doktor iletişimini kolaylaştırmak için kullanılabilir. Örneğin, hastaların sorularını yanıtlayarak, randevu planlamasına yardımcı olarak ve doktorlara hastaların tıbbi geçmişleri hakkında bilgi sağlayarak iletişim süreçlerini daha verimli hale getirebilirler.
- **Tıbbi Teşhis:** BDM'ler, hastaların tıbbi kayıtlarını, semptomlarını ve tıbbi geçmişlerini analiz ederek doktorlara teşhis koymada yardımcı olabilirler.
- **Tedavi Planlaması:** BDM'ler, hastaların tıbbi geçmişlerini ve semptomlarını analiz ederek, kişiselleştirilmiş tedavi planları oluşturabilir ve tedavi süreçlerini optimize edebilirler.

Ancak, BDM'lerin sağlık alanındaki tüm bu potansiyel faydalarını gerçekleştirebilmesi için, bu alana özgü yüksek kaliteli verilerle eğitilmeleri gerekmektedir. Tıbbi metinler, karmaşık terminolojiye, hastalık sınıflandırmalarına, tedavi yöntemlerine ve hasta-doktor iletişim dinamiklerine sahiptir. BDM'lerin bu alandaki verileri doğru bir şekilde anlayabilmesi ve işleyebilmesi için, bu verilere özgü bir şekilde eğitilmeleri gerekmektedir.

Bu noktada, veri toplama ve hazırlama süreçleri büyük önem kazanmaktadır. BDM'lerin sağlık alanına adaptasyonu, bu modellerin sadece genel dil yapısını değil, aynı zamanda sağlık alanına özgü terminolojiyi, hastalık sınıflandırmalarını, tedavi yöntemlerini ve hasta-doktor iletişim dinamiklerini anlamalarını gerektirir. Bu nedenle, BDM'lerin sağlık alanında etkili bir şekilde kullanılabilmesi için, bu alana özgü verilerin toplanması, temizlenmesi, yapılandırılması ve özenle hazırlanması gerekmektedir.

Bu çalışmanın temel amacı, sağlık alanında özelleştirilmiş bir BDM oluşturmak için gerekli veri toplama ve hazırlama süreçlerini detaylı bir şekilde açıklamak ve bu sürecin, BDM'lerin sağlık alanındaki performansını nasıl etkilediğini göstermektir. Bu amaçla iki farklı BDM modeli kullanılacaktır: Meta tarafından geliştirilen LLAMA 3 ve YTÜ COSMOS yapay zeka araştırma grubunun geliştirdiği cosmosGPT v0.1. Bu modellerin farklı yetenekleri ve eğitim verileri, BDM adaptasyon süreci ve sağlık alanına özgü verilerin etkisini daha iyi anlamamızı sağlayacaktır.

LLAMA 3, Meta tarafından geliştirilen açık kaynak kodlu bir BDM'dir. LLAMA 3, geniş bir metin veri kümesi üzerinde eğitilmiş olup, doğal dil işleme görevlerinde etkileyici bir performans sergilemektedir. Ancak, genel amaçlı bir model olması nedeniyle, sağlık alanına özgü terminoloji, hastalık sınıflandırmaları, tedavi yöntemleri ve hasta-doktor iletişim dinamikleri konusunda yeterli bilgiye sahip değildir.

cosmosGPT v0.1 ise, YTÜ COSMOS yapay zeka araştırma grubu tarafından geliştirilmiş ve özellikle Türkçe metinler üzerinde eğitilmiş bir BDM'dir. Bu model, Türkçe dil yapısını ve yaygın kullanılan Türkçe kelime dağarcığını iyi bir şekilde anlamakta ve işlemektedir. Ancak, LLAMA 3 gibi, cosmosGPT v0.1 de sağlık alanına özgü bilgiler konusunda eksikliklere sahiptir.

Bu çalışmada, LLAMA 3 ve cosmosGPT v0.1 modellerinin sağlık alanına adaptasyonu için fine-tuning yöntemi kullanılacaktır. Fine-tuning, önceden eğitilmiş bir BDM'nin, yeni bir göreve veya alana özgü verilerle ek olarak eğitilmesi işlemidir. Bu sayede, model belirli bir alandaki performansını artırabilir. Bu çalışmada kullanılan sağlık alanına özgü veri seti, doktorlar tarafından hastaların sorduğu sorulara verilen cevaplardan oluşmaktadır. Bu veri seti, BDM'lerin sağlık alanındaki terminolojiyi, hastalık sınıflandırmalarını ve hasta-doktor iletişim dinamiklerini öğrenmelerini sağlayacaktır.

Fine-tuning işlemi sırasında, LLAMA 3 ve cosmosGPT v0.1 modelleri bu veri seti üzerinde eğitilecek ve sağlık alanına özgü sorulara daha doğru ve alakalı cevaplar vermesi hedeflenecektir. Modellerin performansı, eğitim ve test veri setleri kullanılarak değerlendirilecektir. Fine-tuning işlemi sonucunda, LLAMA 3 ve cosmosGPT v0.1 modellerinin sağlık alanındaki performansının artması ve sağlık hizmetleri alanında kullanılabilecek özelleştirilmiş BDM'ler olabilmek potansiyelleri gözlemlenecektir.

Başarımın test edilmesi için ayrıca küçük bir veri seti üzerinden modellerin verdiği cevaplar sağlık profesyonelleri tarafından değerlendirilecektir. Bu değerlendirme sonucunda, modellerin sağlık alanındaki terminolojiyi, hastalık sınıflandırmalarını ve hasta-doktor iletişim dinamiklerini ne kadar iyi anladığı ve doğru cevaplar verdiği belirlenecektir.

Bu çalışma, BDM'lerin sağlık hizmetlerinde kullanımı için atılan önemli bir adımdır. BDM'lerin sağlık alanındaki potansiyel faydaları çok çeşitlidir. Ancak, bu potansiyeli tam olarak gerçekleştirebilmek için, doğru ve etkili

veri toplama ve hazırlama süreçlerine dikkat edilmesi gerekmektedir. Bu çalışma, bu süreçleri detaylı bir şekilde açıklayarak, BDM'lerin sağlık alanındaki uygulamalarına yönelik araştırmalara katkı sağlamayı amaçlamaktadır.

2. Veri Toplama ve Hazırlama Süreci:

2.1 Veri Kaynağı ve Toplama: Doktorsitesi.com'dan elde edilen veriler, doktorların anonim profilleri ile hastaların sorduğu sorulara verilen cevaplardan oluşmaktadır. Veri toplama süreci şu adımları içermektedir:

- **Veri Erişimi:** API veya web scraping yöntemleri kullanılarak verilerin toplanması.
- **Veri Güvenliği ve Anonimlik:** Kullanıcıların gizliliğini korumak için tüm kişisel bilgilerin anonimleştirilmesi.
- **Veri Yapılandırması:** Elde edilen verilerin belirli bir formata dönüştürülmesi.

2.2 Veri Temizleme: Toplanan verilerin temizlenmesi, modelin eğitim kalitesini doğrudan etkileyen kritik bir adımdır. Bu adımlar şunları içerir:

- **Boş ve Eksik Değerlerin İşlenmesi:** Boş değerlerin doldurulması veya veri setinden çıkarılması.
- **Hatalı Verilerin Düzeltilmesi:** Yazım hataları ve yanlış bilgilerin düzeltilmesi.
- **Tekrarlayan Verilerin Kaldırılması:** Aynı verilerin birden fazla kez yer almasının önüne geçilmesi.

2.3 Veri Dönüşümü: Verilerin model için uygun hale getirilmesi gerekmektedir. Bu adımlar şunları içerir:

- **Veri Tipi Dönüşümleri:** Verilerin gerekli formatlara dönüştürülmesi (örneğin, tarih formatlarının standart hale getirilmesi).
- **Metin Ön İşleme:** Metin verilerinin temizlenmesi, stop-word'lerin çıkarılması, lemmatizasyon ve tokenizasyon işlemlerinin uygulanması.

2.4 Veri Kalite Kontrolü: Verilerin model eğitime uygun olup olmadığını kontrol etmek için çeşitli kalite kontrol adımları uygulanır:

- **Veri Tutarlılığı:** Verilerin belirli kurallara ve tutarlılığa uygun olup olmadığının kontrol edilmesi.
- **Örnekleme Kontrolü:** Rastgele örneklerin incelenmesi ve kalitelerinin manuel olarak doğrulanması.

3. Model Eğitimi ve Fine-Tuning Süreci:

3.1 Model Seçimi: Bu çalışmada, Meta tarafından geliştirilen LLAMA 3 ve YTÜ COSMOS tarafından geliştirilen cosmosGPT v0.1 modelleri kullanılmaktadır.

3.2 Fine-Tuning Yöntemi: Fine-tuning, önceden eğitilmiş bir modelin, yeni ve spesifik bir görev veya alan için ek eğitim verilerek optimize edilmesi sürecidir. Bu çalışmada, doktorların hastalara verdiği cevaplar ile modellerin sağlık alanına uyum sağlaması amaçlanmıştır.

3.3 Eğitim Süreci: Eğitim süreci şu adımları içermektedir:

- **Eğitim ve Test Veri Setlerinin Oluşturulması:** Verilerin eğitim ve test setlerine ayrılması.
- **Model Parametrelerinin Ayarlanması:** Eğitim sürecinde kullanılacak hiperparametrelerin belirlenmesi.
- **Eğitim:** Modelin eğitim veri seti üzerinde fine-tuning edilmesi.
- **Değerlendirme:** Modelin test veri seti üzerindeki performansının değerlendirilmesi.

4. Sonuçlar ve Değerlendirme:

4.1 Model Performans Değerlendirmesi: Eğitim sonrası modellerin performansı, belirlenen metrikler kullanılarak değerlendirilmiştir:

- **Doğruluk (Accuracy):** Modelin doğru cevap verme oranı.
- **Hassasiyet (Precision) ve Kapsam (Recall):** Modelin verdiği cevapların doğruluğu ve eksiksizliği.
- **F1 Skoru:** Hassasiyet ve kapsamın harmanlanmış ortalaması.

4.2 Sağlık Profesyonellerinin Değerlendirmesi: Modellerin sağlık profesyonelleri tarafından değerlendirilmesi, gerçek dünya uygulamaları için önemli geri bildirimler sağlamıştır. Bu değerlendirme, modellerin sağlık alanındaki terminolojiyi ve iletişim dinamiklerini ne kadar iyi anladığını ortaya koymuştur.

5. Tartışma:

5.1 Elde Edilen Bulguların Analizi: Eğitim sonuçları ve sağlık profesyonellerinin geri bildirimleri analiz edilerek, modellerin performansı ve sağlık alanına adaptasyon süreci değerlendirilmiştir.

5.2 Veri Kalitesinin Önemi: Veri kalitesinin model performansı üzerindeki etkileri tartışılmış ve yüksek kaliteli verilerin önemine vurgu yapılmıştır.

5.3 Gelecekteki Çalışmalar İçin Öneriler: Bu çalışmanın bulguları doğrultusunda, gelecekte yapılacak çalışmalar için öneriler sunulmuştur. Bu öneriler, daha geniş veri setlerinin kullanımı, farklı model mimarilerinin incelenmesi ve modellerin klinik uygulamalarda test edilmesi gibi konuları içermektedir.

6. Sonuç:

BDM'lerin sağlık alanındaki potansiyel faydaları çok çeşitlidir. Bu çalışmada, BDM'lerin sağlık alanına adaptasyonu için veri toplama ve hazırlama süreçleri detaylı bir şekilde açıklanmış ve bu süreçlerin model performansı üzerindeki etkileri değerlendirilmiştir. Sonuçlar, yüksek kaliteli verilerin önemini ve bu verilerin doğru bir şekilde işlenmesi gerektiğini göstermektedir. Bu çalışma, BDM'lerin sağlık hizmetlerinde etkin bir şekilde kullanılabilmesi için atılan önemli bir adımdır ve gelecekteki araştırmalara ışık tutmayı amaçlamaktadır.

Kaynakça:

Burada, çalışmada kullanılan tüm kaynaklar ve referanslar belirtilmelidir. Bu, veri toplama yöntemlerinden, kullanılan algoritmalar ve modeller hakkında teknik dökümanlara kadar geniş bir yelpazeyi kapsamalıdır.