



КЫРГЫЗ - ТҮРК “МАНАС” УНИВЕРСИТЕТИ

ТАБИГЫЙ ИЛИМДЕР ИНСТИТУТУ

КОМПЬЮТЕРДИК ИНЖЕНЕРИЯ БӨЛҮМҮ

ТЕКСТ ТАПШЫРМАЛАРЫН БААЛОО МОДЕЛИН

ИШТЕП ЧЫГУУ

Даярдаган:

Сапаралиева Ырыскүл Едильбековна

Илимий жетекчи:

Доц, Др. Исмаилова Рита

Магистрдик диссертация

Июнь 2022

БИШКЕК, КЫРГЫЗСТАН

ИЛИМИЙ ЭТИКАГА ШАЙКЕШТИГИ

Бул диссертацияда колдонулган бардык маалыматтар академиялык жана этикалык эрежелерге жооп берет. Ошондой эле, башка илимий булактардан алынган маалыматтарга “АДАБИЯТ” бөлүмүндө шилтемелер көрсөтүлдү.

Сапаралиева Ырыскүл

Колу:

BİLİMSEL ETİĞE UYGUNLUK

Bu çalışmadaki tüm bilgilerin, akademik ve etik kurallara uygun bir şekilde elde edildiğini beyan ederim. Aynı zamanda bu kural ve davranışların gerektirdiği gibi, bu çalışmanın özünde olmayan tüm materyal ve sonuçları tam olarak aktardığımı ve referans gösterdiğimi belirtirim.

Ырыскүл САПАРАЛІЕВА

İmza:

ЖОЛ-ЖОБОГО ШАЙКЕШТИГИ

“Текст тапшырмаларын баалоо моделин иштеп чыгуу “ аталышындагы магистрдик диссертация, Кыргыз-Түрк “Манас” университетинин магистрдик диссертация жазуу жол-жобосуна ылайык даярдалды.

Даярдаган:

Сапаралиева Ырыскүл

Колу:

Илимий жетекчи:

Доц., Др. Исмаилова Рита

Колу:

Компьютердик инженерия бөлүм башчысы:

Доц. Др. Райымбек Султанов

Колу:

YÖNERGEYE UYGUNLUK

“ Metin ödevlerin değerlendirme modelinin geliştirilmesi “ adlı Yüksek Lisans Tezi, Kırgızistan-Türkiye Manas Üniversitesi Lisansüstü Tez Önerisi ve Tez Yazım Yönergesi’ne uygun olarak hazırlanmıştır.

Tezi hazırlayan:

İrskül SAPARALIEVA

İmza

Danışman:Doç.Dr.Rita

İSMAİLOVA

İmza

Bilgisayar Mühendisliği ABD Başkanı:

Doç. Dr. Rayımbek SULTANOV

İmza

КАБЫЛ АЛУУ ЖАНА ЧЕЧИМ

Доц. Др. Рита Исмаилованын жетекчилиги астында Ырыскүл Сапаралиева тарабынан даярдалган “Текст тапшырмаларын баалоо моделин иштеп чыгуу” темасындагы диссертация, калыстар тарабынан Кыргыз Түрк “Манас” Университети, Табигый Илимдер Институту Компьютердик инженерия бөлүмүндө Магистрдик диссертация катары кабыл алынды.

..... /..... / 2022

КОМИССИЯ:

Илимий жетекчи Доц. Др. Исмаилова Рита

Төрагасы Доц. Др. Искаков Рысбек

Мүчө Проф. Др. Биримкулов Улан

Мүчө Доц. Др. Шаршембаев Бакыт

Мүчө Доц. Др. Казакбаева Замиргул

Мүчө Ага окут. Др. Жумабаева Чынара

ЧЕЧИМ:

Бул магистрдик иштин кабыл алынышы Институт башкаруу кеңешинин датасына жана санындагы чечими менен бекитилди. /..... / 2022

Доц.Др. Исмет Алтынташ

Институт Мүдүрү

KABUL VE ONAY

Doç.Dr. Rita İsmailova danışmanlığında Iriskül SAPARALIVA tarafından hazırlanan “Metin ödevlerin değerlendirme modelinin geliştirilmesi” adlı bu çalışma, jürimiz tarafından Kırgızistan-Türkiye Manas Üniversitesi Fen Bilimleri Enstitüsü Bilgisayar Mühendisliği AnaBilim Dalında Yüksek Lisans Tezi olarak kabul edilmiştir.

..... /..... / 2022

JÜRİ:

Danışman	Doç.Dr. Rita İSMAİLOVA
Jüri başkanı	Doç.Dr. Rısbek İSKAKOV
Üye	Prof.Dr. Ulan BİRİMKULOV
Üye	Doç.Dr. Bakyt ŞARŞEMBAYEV
Üye	Doç.Dr. Zamırgul KAZAKBAEVA
Üye	Öğr.Gör. Çınara CUMABAEVA

ONAY:

Bu tezin kabulü Enstitü Yönetim Kurulunun tarih ve sayılı kararı ile onaylanmıştır.

..... /..... / 2022

Doç.Dr. İsmet ALTINTAŞ

Enstitü Müdür.

ЫРААЗЫЧЫЛЫК КАТ

Бул иштеги илимий жетекчим Доц. Др. Рита Исмаиловага иш аткаруу этабында көргөзгөн колдоосу, мотивациясы жана туура багыт бергендиги үчүн терең ыраазычылык билдирем. Рита Исмаилованын жетекчилигинин алдында билим деңгээлим жана дүйнө таануум жогорулап, конференцияга катышып, илимий макалалар жарыяланды. Бул жыйынтыктар келечекте дагы көптөгөн илим изилдөөлөрдүн башталышына жана алардын ийгиликтерге жетүүсүнө себеп болот деп ишенем. Компьютердик инженерия бөлүм башчысы Доц.Др. Райымбек Султановго жана ал бөлүмдө эмгектенген баардык эжей-агайларга, чогу окуган группалашым, Нурмила Жетимишовага жана ата-энеме терең урматтоо менен ыраазычылык билдирем.

Июнь, 2022

Сапаралиева Ырыскүл

ТЕКСТ ТАПШЫРМАЛАРЫН БААЛОО МОДЕЛИН ИШТЕП ЧЫГУУ

Сапаралиева Ырыскүл

Кыргыз-Түрк «Манас» университети, Табигый илимдер институту

Магистрдик иш, июнь 2022

Илимий жетекчи: Доц., др. Исмаилова Рита

АННОТАЦИЯ

Чет тилдерди үйрөнүүнүн артыкчылыгы – бул сиздин жашооңузду дээрлик бардык тармагында ийгиликке жетишиңизди шарттайт. Бүгүнкү күнгө чейин, бул, балким, дүйнөдөгү эң пайдалуу чеберчилик. Жаңы тилди үйрөнүү мээңизди жаңы грамматика жана лексика эрежелери менен таанышууга мажбурлайт. Мындан тышкары, чет тилин үйрөнүү жумуш табуу мүмкүнчүлүктөрүн жакшыртат. Эссе (дилбаян) жазуу – жаңы сөздөрдү жаттоого, алардын ортосундагы байланыштарды түзүүгө жана контексттик кырдаалдарда колдонууга эс тутумуңузду үйрөтүүгө мүмкүндүк берген тил үйрөнүүнүн эң маанилүү ыкмаларынын бири. Бирок, педагогдор үчүн бул түр тапшырмаларды текшерүү көп убакытты талап кылат. Экинчи жагынан, текст түрүнө жазылган тапшырмаларды текшерүүнү автоматташтыруу бул процессти бир топ жеңилдетет. Текст тапшырмасын баалоо алгоритминин долбоору - бул табигый тилди иштетүү куралдарынын жардамы менен жазылган дилбаяндар сыяктуу текст тапшырмаларын баалаган система. Бул изилдөөдө табигый тилди иштетүү боюнча изилдөө иштерин өркүндөтүү менен катар, үй тапшырмасын текстке негиздөө менен алгоритм түзүлүп жатат. Студенттер өз эсселерин жаза турган редакторду жана аларды бир нече параметрлердин негизинде баалаган негизги программаны сунуш кылабыз. Теманын актуалдуулугуна көңүл бура турган болсок, ар бир киши бирден көп тил билиши керек экендигин, ааламдаштыруу эрасында учурдун талабы катары

карасак болот. Ал эми тил үйрөнүүнүн бир методологиясы - бул баяндама жазуу. Мугалим тараптан бул процессти карай турган болсок, окуучулардын баяндамасын текшерүү абдан көп убакыт ала турган процесс, жана бул көйгөй, окуучулардын саны көбөйгөй сайын өтө көп убакыт жана мугалимдин көнүлүн талап кылып баштайт. Ал эми акыркы учурда маалымат жана коммуникациянын өнүгүүсүнөн келип чыккан онлайн үйрөнүү жана массалык ачык онлайн курстар (Massive Open Online Courses) аркылуу тил үйрөтүүнү эске ала турган болсок, дилбаяндарды бат текшерүү өтө маанилүү экенин көрө алабыз. Мугалим тарабынан баяндама текшерүүнүн дагы бир көйгөй - инсан фактору, башкача айтканда бул баалоонун субъективдүүлүгү, жана дилбаянга берилген баа мугалимдин баяндамадагы пикирге болгон субъективдүү мамилесинин чагылдырышы мүмкүн.

Муну менен бирге, табигый тил иштетүү - ТТИ (natural language processing - NLP) багытын карай турган босок, бул багытта акыркы жылдарда абдан чоң жетишкендиктер болгонун көрө алабыз. Табигый тилди иштетүүнүн изилдөөчүлөрү, кишилер тилди кантип түшүнүп жана колдонуп атканын изилдөө аркылуу, бул процессте адам баласы колдонгон керектүү иш-аракеттерди кадам кадам түшүнүү аркылуу машиналарга (компьютерлерге) үйрөтүү техникаларды, алгоритмдерди жана аспаптарды иштеп чыгууну көздөмөкчү. Бул процесстин натыйжасында, компьютердик тутумдарды тилге үйрөтүп, керектүү тапшырмаларды түшүнүүгө жана аткарууга мүмкүнчүлүк берүү болуп саналат. Табигый тил иштетүү мисалдары катары машиналык котормо, табигый тилдеги тексттер менен иштөө, көп тилдүү колдонуучунун интерфейстерин даярдоо, сүйлөө таануу, жасалма чалгындоо жана эксперт тутумдары жана башка багыттарды айтып

кетсек болот. Жана албетте, тил үйрөнүүдө табигый тил иштетүү багытынын ролу абдан чоң экенин айтып кете алабыз.

Бул магистрдик диссертация алкагында, окуучулардын баяндама жазуу тапшырмасын баалоо процессин автоматташтыруу, башкача айтканда тапшырма баалоого машина үйрөтүү алгоритмин иштеп чыгуу максаты коюлган. Башкача айтканда, табигый тилди иштетүү куралдарынын жардамы менен окуучулар тарабынан жазылган баяндама же дилбаяны сыяктуу текст тапшырмаларын баалаган системаны иштеп чыгууну жана окуучулар жазган баяндамаларды бир нече параметрлердин негизинде баалаган негизги программаны сунуш кылууну көздөмөкчүбүз. Бул максатка жетүү үчүн, мугалим тарабынан текшерилген баяндамалардын жыйынты топтолуп, баалоо параметрлары иштерип чыкты. Жалпысынан бул процессти автоматташтыруу аракеттери көптөн бери изилдөнүүдө. Азыркы учурда актуалдуу болгон массалык ачык онлайн сабактарда колдонуу мүмкүнчүлүгү менен бирге бул түр долбоорлор көптөгөн изилдөөчүлөрдү кызыктырууда. Бирок, табигый тилди изилдөө жана баалоо үчүн алгоритмдерди иштеп чыгуу абдан оор көйгөйлөрдүн бири болуп саналат. Ошондуктан бул диссертациялык иштин алкагында табигый тил иштетүү негиздери даяр болгон түрк тилиндеги баяндамаларды текшерүүнү автоматташтыруу максаты коюлду.

Бул диссертациялык иште, эссе тапшырмаларын баалоо модели ишке ашырылды. Иш алкагында, табигый тилде жазылган сүйлөмдөрдүн өзгөчөлүктөрүн чыгаруу алгоритмдерин колдонуп, булардын негизинде баалоо жүргүзүлдү. Баалоону божомолдоо үчүн өрнөк тапшырмалардын жыйыны колдонулду жана бул өрнөктөгү баалоо жыйынтыктарынын негизинде машина үйрөтүү методунун

жардамы менен модель курулду. Иштелип чыккан модель, програмдык жабдык же плагин түрүнө капсуляцияланып, Университетибизде колдонулган жабдыктарга интеграцияланышы пландалууда.

Ачык сөздөр: Табигый тил иштетүү, автоматташтырылган баяндама текшерүү, компьютердин жардамы менен тил үйрөнүү, тапшырма баалоо тутумдары



METİN ÖDEVLERİN DEĞERLENDİRME MODELİNİN GELİŞTİRİLMESİ

İrışköl SAPARALIEVA

Kirgizistan-Türkiye Manas Üniversitesi, Fen Bilimleri Enstitüsü

Yüksek Lisans Tezi, Haziran 2022

Danışman: Doç. Dr. Rita İSMAİLOVA

GENİŞ ÖZET

Yabancı dil öğrenmenin avantajı, hayatınızın hemen her alanında sizi başarıya hazırlayabilmesidir. Bugüne kadar, bu belki de şimdiye kadarki en faydalı gerçek dünya becerisidir. Yeni bir dil öğrenmek, beyninizi yeni gramer ve kelime kurallarına aşina olmaya zorlar. Ek olarak, yabancı dil öğrenmek iş beklentilerinizi iyileştirebilir. Kompozisyon yazma, yeni kelimeleri ezberlemek, aralarında bağlantı kurmak ve bunları bağlamsal durumlarda kullanmak için hafızanızı eğitmenize olanak tanıyan en önemli dil öğrenme yöntemlerinden biridir. Bununla birlikte, eğitimciler için, öğrenci kompozisyonlarının değerlendirilmesi zaman alan bir süreçtir. Öte yandan, öğrenci kompozisyonlarının yada başka metin yazma ödevlerin değerlendirilmesini otomatikleştirmek bu süreci büyük ölçüde kolaylaştırabilir. Bir metnin dilbilgisin değerlendirme algoritması, doğal dil işleme araçlarının yardımıyla yazılan kompozisyon gibi metin yada makaleleri değerlendiren bir sistemdir. Bu çalışmada, doğal dil kullanımına yönelik araştırmaların geliştirilmesinin yanı sıra, öğrenci kompozisyonlarının yola çıkılarak bir algoritma geliştirilmektedir. İlk adım olarak öğrencilerin kompozisyonlarının yazmaları için bir editör ve bunları çeşitli parametrelere göre değerlendiren temel bir program sunuyoruz.

Konunun aciliyeti göz önüne alındığında herkesin birden fazla dil bilmesi gerektiği küreselleşme çağında çağın bir gereği olarak değerlendirilebilir. Ve dil öğrenmenin bir yöntemi de rapor yazmaktır. Bu sürece öğretmen gözüyle bakacak olursak, öğrenci

raporlarını kontrol etmek oldukça zaman alan bir işlemdir ve bu sorun öğrenci sayısı arttıkça çok fazla zaman ve öğretmenin dikkatini gerektirmeye başlar. Ve kitlesel Açık Çevrimiçi Kurslar aracılığıyla bilgi ve iletişim, çevrimiçi öğrenme ve dil öğretimindeki son gelişmeleri dikkate alırsak, makaleleri hızlı bir şekilde kontrol etmenin çok önemli olduğunu görebiliriz. Raporun öğretmen incelemesinin bir başka sorunu da kişilik faktörüdür, yani bu değerlendirmenin öznelliği ve makalenin değerlendirilmesi, öğretmenin rapordaki görüşe yönelik öznel tutumunu yansıtabilir.

Aynı zamanda, doğal dil işleme (NLP) için eşik eşiktir ve son yıllarda bu alanda büyük başarılar elde edildiğini görebiliriz. Doğal dil işleme araştırmacıları, insanların dili nasıl anladığını ve kullandığını inceleyerek ve insanların bu süreçte kullandığı adımları adım adım anlayarak makinelere (bilgisayarlara) öğretmek için teknikler, algoritmalar ve araçlar geliştirmeyi amaçlar. Bu sürecin sonucu, bilgisayar sistemlerine dili öğretmek ve gerekli görevleri anlamalarını ve gerçekleştirmelerini sağlamaktır. Doğal dil işleme örnekleri arasında makine çevirisi, doğal dil metinleriyle çalışma, çok dilli kullanıcı arabirimleri geliştirme, konuşma tanıma, yapay zeka ve uzman sistemler ve daha fazlası yer alır. Ve tabii ki doğal dil işlemenin dil öğrenimindeki rolünün çok önemli olduğunu söyleyebiliriz.

Bu yüksek lisans tezinin amacı, öğrencilerin rapor yazma görevlerini değerlendirme sürecini otomatikleştirmek, yani görev değerlendirmesi için makine öğrenimi için bir algoritma geliştirmektir. Başka bir deyişle, öğrenci raporları veya kompozisyonlar gibi metin ödevlerini doğal dil işleme araçlarını kullanarak değerlendirmek için bir sistem geliştirmeyi ve öğrenci raporlarını çeşitli parametrelere göre değerlendiren temel bir program sunmayı amaçlıyoruz. Bu amaca ulaşmak için öğretmen tarafından gözden geçirilen raporlar toplandı ve değerlendirme parametreleri geliştirildi. Genel olarak, bu

süreci otomatikleştirme girişimleri uzun süredir incelenmiştir. Birçok araştırmacı, bu tür projelerle ve bunları şu anda ilgili olan toplu açık çevrimiçi sınıflarda kullanma olasılığıyla ilgilenmektedir. Bununla birlikte, doğal dilin incelenmesi ve değerlendirilmesi için algoritmaların geliştirilmesi en zor zorluklardan biridir. Bu nedenle, bu tezin amacı, doğal dil işleme temellerine sahip Türkçe sunuların incelenmesini otomatik hale getirmektir.

Bu tez çalışmasında kompozisyon ödevlerini değerlendirmek için bir model uygulanmıştır. Çalışma kapsamında, doğal dilde yazılan cümlelerin özelliklerini çıkarmak için algoritmaların kullanımına dayalı bir değerlendirme yapılmıştır. Değerlendirmeyi tahmin etmek için bir dizi örnek görev kullanılmış ve bu modeldeki değerlendirmenin sonuçlarına dayalı olarak makine öğrenmesi yöntemi kullanılarak bir model oluşturulmuştur. Geliştirilen modelin yazılım veya eklenti şeklinde kapsüllenmesi ve Üniversitemizde kullanılan donanımlara entegre edilmesi planlanmaktadır.

Anahtar Kelimeler: Doğal dil işleme, otomatik konuşma testi, bilgisayar destekli dil öğrenimi, görev değerlendirme sistemleri, makine öğrenimi, deneme doğrulama algoritması geliştirme projesi

DEVELOPMENT OF TEXT ASSIGNMENTS EVALUATION MODEL

Yryskul SAPARALIEVA

**Kyrgyz-Turkish Manas university, Graduate school of Natural
and Applied Sciences**

M.S.Thesis, June 2022

Supervisor: Dr. Prof. Rita ISMAILOVA

ABSTRACT

We all know that learning foreign languages can set students up for success in almost every area of your life. Knowing many languages is perhaps the most useful real-world skill ever. Learning a new language forces your brain to become familiar with new grammar and vocabulary rules. In addition, learning a foreign language can improve your job prospects. Essay writing is one of the most important language learning methods that allows you to train your memory to memorize new words, make connections between them and use them in contextual situations. However, for educators, proofreading and evaluating essays is a time-consuming process. On the other hand, automatization of an essay grading process could greatly facilitate this process. An essay grading algorithm is a system that grades text tasks, such as essays, written using natural language processing tools. In this study, in addition to improving research on the use of natural language, an essay grading algorithm is developed.

The study on essay grading is hard to underestimate; the fact that everyone should know more than one language can be considered a requirement of the times in the era of globalization. And one of the methods of learning a language is writing essays. As mentioned above, from a teacher's perspective, grading students' essays is a very time-consuming process, and this problem begins to require a lot of time and attention from the teacher as the number of students increases. And if we take into account the recent

development of information and communication technologies, online learning, as well as teaching languages through massive open online courses, we can see that it is very important to grade essays quickly. Another problem with teachers' grading of the students' essays is a co-called human factor, i.e. the subjectivity of this assessment as the assessment of an essay may reflect the subjective attitude of the teacher to the opinion reflected in an essay rather than its grammar.

On the other hand, natural language processing (NLP), which is defined in literature as “the branch of computer science concerned with giving computers the ability to understand text and spoken words in much the same way human beings can”, has shown a profound development in recent years. Natural language processing researchers aim to develop methods, algorithms, and tools for learning machines (computers) by studying how humans understand and use language and understanding the steps humans use in the process. The result of this process is to “teach” computer systems the language and enable them to understand and perform the required tasks. Examples of natural language processing applications include machine translation, natural language processing, multilingual user interface development, speech recognition, artificial intelligence and expert systems, and more. And, of course, we can say that the role of natural language processing in language learning is very important.

The purpose of this master's thesis is to automate the process of grading students' essay assignments, i.e. development of a machine learning algorithm for essay grading. In other words, we intend to develop a system for grading essays or any other textual assignments using natural language processing tools, and offer a basic program that grades students' essays based on several grammar parameters. To achieve this goal, teacher-reviewed essays were collected and assessment parameters for machine grading were developed.

In general, attempts to automate this process have long been worked out. Many researchers are interested in these types of projects, as well as the possibility of their use in massive open online classes, which is relevant at the present time of MOOCs. However, the development of algorithms for learning and evaluating natural language is one of the most difficult tasks. Therefore, the aim of this dissertation is to automate the grading of essays of students who learn Turkish.

In this dissertation work, a model for grading essay assignments was implemented. As part of the work, an assessment was carried out based on the use of algorithms for extracting features of sentences written in natural language. A set of typical tasks was used to predict the assessment and based on the results of the assessment, a machine learning model was built in this model. The developed model was encapsulated in the form of software or a plug-in and integrated into the equipment used at our university.

Key words: natural language processing, automated speech testing, computer assisted language teaching, text assignments evaluating system.

РАЗРАБОТКА МОДЕЛИ ОЦЕНКИ ТЕКСТОВЫХ ЗАДАНИЙ

Сапаралиева Ырыскүл

Кыргызко-Турецкий университет «Манас»,

Институт Естественных наук

Магистерская диссертация, июнь 2022

Научный руководитель: Др. Доц. Рита Исмаилова

АННОТАЦИЯ

В соответствии с последними тенденциями в Интернете вещей компьютеры и другие устройства программируются так, чтобы общаться, слушать и говорить, как люди. Таким образом, область обработки естественного языка становится очень важной для изучения этого типа искусственного интеллекта. Алгоритм оценки текстового задания — это система, которая оценивает текстовые задания, например эссе, написанные с помощью инструментов обработки естественного языка. В данном исследовании, помимо совершенствования исследований по использованию естественного языка, разрабатывается алгоритм на основе текста домашнего задания. Мы предлагаем редактор для студентов, чтобы написать свои эссе и базовую программу, которая оценивает их по нескольким параметрам.

Учитывая актуальность темы, то, что каждый должен знать более одного языка, можно считать требованием времени в эпоху глобализации. И одним из методов изучения языка является написание доклада. Если посмотреть на этот процесс с точки зрения преподавателя, то проверка отчетов студентов является очень трудоемким процессом, и эта проблема начинает требовать много времени и внимания преподавателя по мере увеличения количества студентов. И если мы примем во внимание недавнее развитие информации и коммуникации, онлайн-обучения и преподавания языков через массовые открытые онлайн-курсы, мы увидим, что очень важно быстро проверять эссе. Еще одной проблемой

рецензирования учителем отчета является личностный фактор, т.е. субъективность этой оценки, а оценка эссе может отражать субъективное отношение учителя к мнению в отчете.

В то же время порогом обработки естественного языка (ОБЯ) является порог, и мы видим, что за последние годы в этой области были достигнуты большие успехи. Исследователи обработки естественного языка стремятся разработать методы, алгоритмы и инструменты для обучения машин (компьютеров), изучая, как люди понимают и используют язык, и шаг за шагом понимая шаги, которые люди используют в этом процессе. Результатом этого процесса является обучение компьютерных систем языку и предоставление им возможности понимать и выполнять требуемые задачи. Примеры обработки естественного языка включают машинный перевод, работу с текстами на естественном языке, разработку многоязычных пользовательских интерфейсов, распознавание речи, искусственный интеллект и экспертные системы и многое другое. И, конечно же, мы можем сказать, что роль обработки естественного языка в изучении языка очень важна.

Целью данной магистерской диссертации является автоматизация процесса оценивания заданий студентов по написанию отчетов, т.е. разработка алгоритма машинного обучения для оценивания заданий. Другими словами, мы намерены разработать систему оценивания текстовых заданий, таких как отчеты учащихся или рефераты, с использованием средств обработки естественного языка, и предложить базовую программу, оценивающую отчеты учащихся по нескольким параметрам. Для достижения этой цели были собраны эссе, проверенные учителями, и разработаны параметры оценки. Вообще попытки автоматизировать

этот процесс давно прорабатываются. Многих исследователей интересуют такие виды проектов, а также возможность их использования в массовых открытых онлайн-занятиях, что актуально в настоящее время. Однако разработка алгоритмов для изучения и оценки естественного языка является одной из самых сложных задач. Поэтому целью данной диссертации является автоматизация проверки презентаций на турецком языке, имеющих основы обработки естественного языка. В данной диссертационной работе была реализована модель оценивания заданий эссе. В рамках работы была проведена оценка, основанная на использовании алгоритмов для извлечения признаков предложений, написанных на естественном языке. Для прогнозирования оценки использовался набор типовых задач, и по результатам оценки в этой модели была построена модель методом машинного обучения. Разработанную модель планируется инкапсулировать в виде программного обеспечения или плагина и интегрировать в оборудование, используемое в нашем университете.

Ключевые слова: обработка естественного языка, автоматизированное речевое тестирование, компьютерное обучение языку, системы оценки текстовых заданий.

МАЗМУНУ

ТЕКСТ ТАПШЫРМАЛАРЫН БААЛОО МОДЕЛИН

ИШТЕП ЧЫГУУ

ИЛИМИЙ ЭТИКАГА ШАЙКЕШТИГИ.....	ii
BİLİMSEL ETİĞE UYGUNLUK.....	ii
ЖОЛ-ЖОБОГО ШАЙКЕШТИГИ.....	iii
YÖNERGEYE UYGUNLUK.....	ivi
КАБЫЛ АЛУУ ЖАНА ЧЕЧИМ.....	iv
KABUL VE ONAY.....	vi
ЫРААЗЫЧЫЛЫК КАТ.....	vi
АННОТАЦИЯ.....	vii
GENİŞ ÖZET.....	xi
ABSTRACT.....	xiv
АННОТАЦИЯ.....	xvii
МАЗМУНУ.....	xx
ЖАДЫБАЛДАРДЫН ТИЗМЕСИ.....	xxiii
СҮРӨТТҮН ТИЗМЕСИ.....	xxiv
КЫСКАРТМАЛАР ТИЗМЕСИ.....	xxv
КИРИШҮҮ.....	1

БӨЛҮМ I.

ТАБИГЫЙ ТИЛ ИШТЕТҮҮ ТУУРАЛУУ ЖАЛПЫ МААЛЫМАТ.....	5
1.1. Табигый тилди иштетүү (natural language processing - NLP).....	5
1.2. Табигый тилди компьютердин жардамы менен түшүнүү.....	6
1.3. Табигый тилди иштетүү куралдары жана ыкмалары.....	6

БӨЛҮМ II.

МАШИНАЛЫК ҮЙРӨТҮҮ ТУУРАЛУУ ЖАЛПЫ МААЛЫМАТ	8
2.1. Изилдөөдө үчүн машинаны үйрөтүүнү тандоо себептери	8
2.2. Машиналык үйрөтүү жөнүндө түшүнүк	9
2.3. Машиналык үйрөтүүнүн түрлөрү	10
2.4. Машиналык үйрөтүүнүн алгоритмдери	13

БӨЛҮМ III.

АДАБИЯТ ТАЛДОО	16
3.1. Баяндаманы текшерүүнү автоматташтыруу тарыхы	16
3.2. Баяндаманы текшерүүнү автоматташтырууда колдонулган машина үйрөтүү алгоритмдери	17

БӨЛҮМ IV.

МЕТОД ЖАНА МАТЕРИАЛДАР	21
4.1. Метод жана материалдар	21
4.1.1. Тил китепканалары. Zemberek китепканасы	21
4.1.2. Колдонуучу интерфейсти иштеп чыгуу	23
4.1.3. Машина үйрөтүүдө колдонулган технологиялар	26
4.2. Изилдөө кадамдары	27

БӨЛҮМ V.

БЕРИЛИШТЕРДИ ДАЯРДОО ЖАНА МАШИНАЛЫК ҮЙРӨТҮҮ	29
5.1. Берилиштерди даярдоо	29
5.2. Машиналык үйрөтүү процесси	32

БӨЛҮМ VI.

ЖЫЙЫНТЫКТАР	38
-------------------	----

6.1. Машина үйрөтүү процессинде алынган алгач жыйынтыктар.....	38
6.2. Баалоо тактыктары	40

БӨЛҮМ VII.

ТАЛКУЛООЛОР ЖАНА КОРУТУНДУ	43
7.1 Талкулоолор.....	43
7.2. Корутунду	43
7.3. Алдыдагы изилдөөлөр	44
КОЛДОНУЛГАН БУЛАКТАРДЫН ТИЗМЕСИ:	45
ӨМҮР БАЯН.....	Hata! Yer işareti tanımlanmamış.

ЖАДЫБАЛДАРДЫН ТИЗМЕСИ

Жадыбал 1. Машиналык үйрөтүүгө керек болгон дилбаяндардын көптүгү.....	40
Жадыбал 2. Баалоо тактыктары.....	41
Жадыбал 3. QWK – моделинен алынган жыйынтыктар.....	42



СҮРӨТТҮН ТИЗМЕСИ

Сүрөт 1. Мугалимдин жардамы менен машиналык үйрөтүү схемасы.	12
Сүрөт 2. IG (өтпөй калды), G (өттү), VG(артыкчылык менен өттү) жана MVG (эң жакшы)	20
Сүрөт 3. Ишке ашырылган программдык жабдыктын интерфейси.....	26
Сүрөт 4. xml форматына өткөрүлгөн дилбаяндын көрүнүшү	30
Сүрөт 5. .xml форматына өткөрүлгөн дилбаяндардын массиви.....	31
Сүрөт 6. .xlsx форматында саталган жыйынтыктар	32
Сүрөт 7. Нейрондук тармакты ишке ашыруу схемасы	33
Сүрөт 8. Жоготуулардын эпохадан көз карандылыгы	35
Сүрөт 9. Нейрондук тармактын модели.....	35
Сүрөт 10. Дилбаянды автоматтык түрдө текшерүү системасы.....	36
Сүрөт 11. Кошумча ПЖ чыгарган жыйынтыктардын массиви.....	39
Сүрөт 12. Кошумча ПЖ чыгарган жыйынтыктардын массивинин уландысы	39

КЫСКАРТМАЛАР ТИЗМЕСИ

Кыскартылышы	Чечмелениши
ТТИ	Табигый тилди иштетүү
NLP	Natural language processing
NAACL	Түндүк Америка эсептөөчү лингвистика ассоциациясынын бөлүмү
EACL	Европалык Эсептөөчү лингвистика ассоциациясынын бөлүмү
ML	Machine Learning
AI, ЖИ	Artificial intelligence Жасалма интеллект
PEG	Project Essay Grading
ETS	Educational Testing Service
QWK	Quadratic-weighted kappa
LSA	Latent Semantic Analysis
PLSA	Probabilistic Latent Semantic Analysis
ПЖ	Программдык жабдык

КИРИШҮҮ

Ар бир киши бирден көп тил билиши керек экендигин, ааламдаштыруу эрасында учурдун талабы катары карасак болот. Ал эми тил үйрөнүүнүн бир методологиясы - бул баяндама жазуу. Мугалим тараптан бул процессти карай турган болсок, окуучулардын баяндамасын текшерүү абдан көп убакыт ала турган процесс, жана бул көйгөй, окуучулардын саны көбөйгөй сайын өтө көп убакыт жана мугалимдин көңүлүн талап кылып баштайт. Акыркы учурда маалымат жана коммуникациянын өнүгүүсүнөн келип чыккан онлайн үйрөнүү жана массалык ачык онлайн курстар (Massive Open Online Courses) аркылуу тил үйрөтүүдөн келип чыгат. Бул онлайн сабактардын популярдуулугун, Интернет аркылуу ийкемдүү күн тартиби менен жана көптөгөн учурда акысыз болгондугу менен түшүндүрсөк болот. Бул түр онлайн сабактар, Интернет аркылуу көптөгөн сандагы окуучулар үчүн иштелип чыккан, жана бул сабактарга катышуунун жалгыз шарты - интернет байланышы болгону айтылып келет [1]. Ал эми тил үйрөтүү массалык ачык онлайн курстар (же LMOOC - Language Massive Open Online Courses) - бул чексиз катышуучуларга чексиз катышуу мүмкүнчүлүгүн берген экинчи тил үйрөнүүгө арналган Интернеттеги онлайн курстар [1]. Бул курстар аркылуу көптөгөн дисциплинаны үйрөнүү мүмкүнчүлүгү болгон менен, бир нече авторлор тил үйрөтүүдө баяндама текшерүү өз алдынча да кыйынчылыктарга туш болгондуктан, Интернеттеги айлана-чөйрөдө массалык түргө айлангандыгынан бул чоң көйгөйгө айланаарын жөнүндө айтылууда [3].

Мугалим тарабынан баяндама текшерүүнүн дагы бир көйгөй - инсан фактору, башкача айтканда бул баалоонун субъективдүүлүгү, анткени инсан катары мугалим жалгыз гана тил билүү жөндөмү эмес, ошондой эле учурда баяндамада

келтирилген окуучунун көз карашын кошо баалайт жана бул окуучунун алган баасына таасирин тийгизиши мүмкүн [4]. Башкача айтканда, баалоо студенттин тил билими эмес, мугалимдин баяндамадагы пикирге болгон субъективдүү мамилесинин чагылдыруусу болушу мүмкүн.

Муну менен бирге, табигый тил иштетүү - ТТИ (natural language processing - NLP) багытын карай турган босок, бул багытта акыркы жылдарда абдан чоң жетишкендиктер болгонун көрө алабыз. Табигый тилдерди иштеп чыгуу адам тилин үйрөнүү, түшүнүү жана өндүрүү максатында эсептөө техникасын колдонот [5]. Табигый тилди иштетүүнүн изилдөөчүлөрү, кишилер тилди кантип түшүнүп жана колдонуп атканын изилдөө аркылуу, бул процессте адам баласы колдонгон керектүү иш-аракеттерди кадам кадам түшүнүү аркылуу машиналарга (компьютерлерге) үйрөтүү техникаларды, алгоритмдерди жана аспаптарды иштеп чыгууну көздөмөкчү. Бул процесстин натыйжасында, компьютердик тутумдарды тилге үйрөтүп, керектүү тапшырмаларды түшүнүүгө жана аткарууга мүмкүнчүлүк берүү болуп саналат.

Табигый тил иштетүүнүн пайдубалдары компьютер жана маалымат илимдери, лингвистика, математика, электрдик жана электрондук инженердик, жасалма инженерия жана робототехника, психология ж.б. бир катар дисциплиналар түзөт [6]. Албетте, табигый тил өтө маанилүү болгондуктан, бул компьютер тутумдарда табигый тил иштетүү дагы көптөгөн колдонмо багыттарында да колдонулуп келмекчи. Мисал катары машиналык котормо, табигый тилдеги тексттер менен иштөө, көп тилдүү колдонуучунун интерфейстерин даярдоо, сүйлөө таануу, жасалма чалгындоо жана эксперт тутумдары жана башка багыттарды айтып кетсек

болот. Жана албетте, тил үйрөнүүдө табигый тил иштетүү багытынын ролу абдан чоң экенин айтып кете алабыз.

Бул магистрдик диссертация алкагында, окуучулардын баяндама жазуу тапшырмасын баалоо процессин автоматташтыруу, башкача айтканда тапшырма баалоого машина үйрөтүү алгоритмин иштеп чыгуу максаты коюлган. Башкача айтканда, табигый тилди иштетүү куралдарынын жардамы менен окуучулар тарабынан жазылган баяндама же дилбаяны сыяктуу текст тапшырмаларын баалаган системаны иштеп чыгууну жана окуучулар жазган баяндамаларды бир нече параметрлердин негизинде баалаган негизги программаны сунуш кылууну көздөмөкчүбүз. Бул максатка жетүү үчүн, мугалим тарабынан текшерилген баяндамалардын жыйынты топтолуп, баалоо параметрлары иштерип чыкты. Жалпысынан бул процессти автоматташтыруу аракеттери көптөн бери изилдөнүүдө. Азыркы учурда актуалдуу болгон массалык ачык онлайн сабактарда колдонуу мүмкүнчүлүгү менен бирге бул түр долбоорлор көптөгөн изилдөөчүлөрдү кызыктырууда. Бирок, табигый тилди изилдөө жана баалоо үчүн алгоритмдерди иштеп чыгуу абдан оор көйгөйлөрдүн бири болуп саналат. Ошондуктан бул диссертациялык иштин алкагында табигый тил иштетүү негиздери даяр болгон түрк тилиндеги баяндамаларды текшерүүнү автоматташтыруу максаты коюлду.

Бул магистрдик диссертациянын түзүлүшүнө көңүл бура турган болсок, биринчи бөлүмдө табигый тил иштетүү тууралуу жалпы маалымат берилмекчи. Экинчи бөлүмдө, машинаны үйрөтүү ыкмалары талданып, анын түрлөрү жана негизги колдонулган алгоритдери каралып чыкмакчы. Азыркы учурга чейин баяндамаларды текшерүүнү автоматташтыруу аракеттери, изилдөөлөр жана

баяндама текшерүү үчүн иштелип чыккан програмдык жабдыктар, алардын күчтүү жана күчсүз тараптары үчүнчү бөлүмдө талкууланган. Төртүнчү бөлүмдө, түрк тилин үйрөнүү үчүн жазылган баяндамаларды текшерүү үчүн иштелип чыккан машина үйрөтүү алгоритмдин түзүүнүн жалпы кадамдары жана методологиясы көрсөтүлүп, ал эми бешинчи бөлүмдө, машина үйрөтүү үчүн берилиштерди даярдоо менен бирге машина үйрөтүү процесси тууралуу маалымат берилген. Иштелип чыккан аспаптын иштөө жыйынтыгы алтынчы бөлүмдө көрсөтүлүп, алтынчы бөлүм, магистрдик диссертациянын алкагында иштелип чыккан аспап боюнча талкуулолор жана жаңы илимий изилдөө багыттары боюнча сунуштарыбызды камтымакчы.

БӨЛҮМ I.

ТАБИГЫЙ ТИЛ ИШТЕТҮҮ ТУУРАЛУУ ЖАЛПЫ МААЛЫМАТ

1.1. Табигый тилди иштетүү (natural language processing - NLP)

Табигый Тилди Иштетүү (ТТИ) – бул компьютерлер табигый тилде жазылган текстти же кепти түшүнүү жана иштетүүдө, тапшырмаларды кантип аткара турганын изилдеген изилдөө жана колдонуу чөйрөсү. NLP изилдөөчүлөрү адамдардын тилди кантип түшүнүп, колдонгону жөнүндө маалымат чогултууга умтулушат, ошондуктан компьютер тутумдарына табигый тилдерди түшүндүрүп, каалаган тапшырмаларды аткаруу үчүн манипуляциялоо мүмкүнчүлүгүн берүү менен тийиштүү инструменттер жана ыкмалар иштелип чыгат. NLP негиздери бир катар дисциплиналарда, башкача айтканда компьютер жана маалымат технологиялары, лингвистика, математика, электр жана электроника инженериясы, жасалма интеллект жана робототехника, психология ж.б. издөө маалыматы (CLIR), кеп же речь таануу, жасалма интеллект жана эксперттик системалар жана башка аймактарда колдонулат.

NLP ни колдонуу дүйнөлүк желенин жана китепканалардын жайылышынан улам бир топ белгилүү болуп калды. Бир нече изилдөөчүлөр интернеттин жана китепканалардын бардык артыкчылыктарын пайдалануу үчүн көп тилдүү же кайчылаш тилдүү маалыматты издөөнү, анын ичинде көп тилдүү текстти иштетүү жана көп тилдүү колдонуучу интерфейс системаларын жеңилдетүү үчүн тийиштүү изилдөөлөрдүн зарылдыгын белгилешти.

Ар кандай NLP тапшырмасынын максатында табигый тилди түшүнүү маселеси турат. Табигый тилди түшүнгөн компьютердик программаларды түзүү процесси үч негизги этапты камтыйт: ой жүгүртүү процессин, тил киргизүүнүн өзгөчөлүгүн жана маанисин, дүйнө таанууну. Ошентип, NLP системасы сөздүк деңгээлден

баштай алат - сөздүн морфологиялык түзүлүшүн, мүнөзүн (мисалы, сөз түркүмүн, маанисин) аныктайт, андан кийин сүйлөм деңгээлине өтө алат. Сүйлөмдүн курамындагы ар бир сөздүн сөз түркүмүн, сүйлөмдүн грамматикалык түзүлүшүн, ар бир сүйлөмдүн маанисин аныктайт, андан кийин контекст жана жалпы маңызын же маанисин түшүнөт. Берилген сөз, сүйлөм контекстте белгилүү бир мааниге же болбосо коннотацияга ээ болушу мүмкүн жана башка көптөгөн сөздөр, сүйлөмдөр менен байланыштырылышы мүмкүн.

1.2. Табигый тилди компьютердин жардамы менен түшүнүү

Liddy (1998) жана Feldman (1999) адамдар табигый тилди түшүнүү үчүн колдонгон төмөнкү жети өз ара көз каранды деңгээлдерди ажырата билүү маанилүү деп эсептешет:

- фонетикалык же фонологиялык деңгээл, тыбыштардын туура айтылуусун көзөмөлдөйт;
- морфологиялык деңгээл, суффикс жана префикс, сөздөрдүн маанисин туюнткан эң майда бөлүктөр;
- лексикалык деңгээл, сөздөрдүн жана сүйлөм мүчөлөрүнүн лексикалык маанисин талдайт;
- синтаксистик деңгээл, сүйлөмдүн грамматикасын жана сүйлөмдүн түзүлүшүн аныктайт;
- семантикалык деңгээл, сөздөрдүн жана сүйлөмдөрдүн мааниси боюнча байланышкандыгын тактайт;
- дискурстук деңгээл, документти колдонуу менен тексттин ар кандай түрлөрүнүн структурасын карайт;
- прагматикалык деңгээл, адабий тилге оозеки тилден кирген маалыматтардын жыйынын камтыйт.

1.3. Табигый тилди иштетүү куралдары жана ыкмалары

Бир катар изилдөөчүлөр NLP ишинин маанилүү бөлүгү болгон ар кандай иш-чараларды аткаруу үчүн жакшыртылган технологияны ойлоп табууга аракет кылышкан. Бул эмгектерди төмөнкүдөй классификациялоого болот:

- Лексикалык жана морфологиялык талдоо, сөз айкаштарын түзүү, сөздөрдү сегменттөө ж.б;
- Семантикалык жана дискурсивдүү анализ, сөздүн мааниси жана билимди чагылдыруу;
- Билимге негизделген NLP ыкмалары жана куралдары

NLP эксперименталдык системалар өтө аз учурда реалдуу системаларга же продуктыларга айландырылат. NLP изилдөөлөрүндөгү негизги көйгөйлөрдүн бири, маалымат издөө сыяктуу тармактарда, чоң сыноо топтомдорунун жана көп жолу колдонула турган эксперименталдык методдордун жана куралдардын жоктугу болууда. Бирок, акыркы бир нече жылда кырдаал өзгөрдү. Учурда бир нече улуттук жана эл аралык изилдөө топтору тесттердин чоң коллекцияларын, ошондой эле эксперименталдык инструменттерди жана методдорду түзүү жана кайра колдонуу үчүн биргелешип иштеп жатышат. Маалыматты бөлүшүү конференциялары пайда болгондон бери, TREC топтук изилдөө аракеттеринин сериялары жана NAACL (Түндүк Америка эсептөөчү лингвистика ассоциациясынын бөлүмү), EACL (Европалык эсептөөчү лингвистика ассоциациясынын бөлүмү) үзгүлтүксүз конференциялар жана семинарлар менен кеңейди. Бул топтук изилдөөлөр изилдөөчүлөргө көп жолу колдонулуучу NLP куралдарын, тест топтомун жана эксперименталдык методологияларды түзүү аркылуу өз тажрыйбалары менен бөлүшүүгө жардам берет [7].

БӨЛҮМ II.

МАШИНАЛЫК ҮЙРӨТҮҮ ТУУРАЛУУ ЖАЛПЫ МААЛЫМАТ

2.1. Изилдөөдө үчүн машинаны үйрөтүүнү тандоо себептери

Тил үйрөнүүдө, эссе жазуу негизги тапшырмалардын түрү. Бирок бул түр тапшырмаларды текшерүү абдан көп убакыт алуучу процесс. Ошондуктан бул процессти автоматташтыруу аракеттери көптөн бери изилдөнүүдө. Азыркы учурда актуалдуу болгон массалык ачык онлайн сабактарда колдонуу мүмкүнчүлүгү менен бирге бул түр долбоорлор көптөгөн изилдөөчүлөрдү кызыктырууда. Бирок, табигый тилди изилдөө жана баалоо үчүн алгоритмдерди иштеп чыгуу абдан оор көйгөйлөрдүн бири болуп саналат. Эссе текшерүү алгоритмин иштеп чыгуу проекти, бул проблеманы чечүү үчүн өнүктүрүлгөн, дилбаяндарды автоматтык түрдө текшере алган система. Бул система азырынча Түркчө жазылган дилбаяндар үчүн гана иштөөдө. Бул долбоорду ишке ашырууда машиналык үйрөтүү алгоритмдери колдонулду. Machine Learning (ML) - бул компьютердик системалар ачык программаланбастан, белгилүү бир тапшырманы аткаруу үчүн колдонгон алгоритмдерди жана статистикалык моделдердин үстүнөн жүргүзүлгөн илимий изилдөөлөр. Биз күнүмдүк колдонгон көптөгөн тиркемелердеги алгоритмдер машиналык үйрөтүүнү колдонушат. Google сыяктуу издөө системасы интернеттен маалымат издөө үчүн бул түрдүү алгоритмдерди колдонот жана анын жакшы иштешинин себептеринин бири - бул машиналык үйрөтүү алгоритминин веб-баракчаларынын рейтингдерин туура үйрөнгөндүгүндө. Бул алгоритмдер ар кандай максаттар үчүн колдонулат, мисалы, маалыматтарды иштетүү, сүрөттөрдү иштетүү, болжолдоочу аналитика ж.б. Машина үйрөнүүнү колдонуунун негизги артыкчылыгы - алгоритм маалыматтар менен эмне кылуу керектигин билгенден кийин, ал өз ишин автоматтык түрдө аткара алат [8].

2.2. Машиналык үйрөтүү жөнүндө түшүнүк

Жасалма интеллект (AI) - адамдын интеллекттин өрнөк катры кабыл алып, адамдын интеллектеги жасай алган иш аракеттерди программдык жабдыктын жардамы аркылуу ишке ашыра алган илимдин бир багыты. Жасалма интеллектти ишке ашыруу үчүн ар түрдүү методдор ишке ашырылган. Жасалма интеллект Data science жана Machine learning методдорунун бардык компоненттерин камтыйт.

Машиналык үйрөтүү - бул жасалма интеллекттин бөлүмдөрүнүн бири. Машиналык үйрөтүү так аныкталган эрежелерди сактабастан, берилген маалыматтарга негизделген жыйынтык чыгарууга мүмкүндүк берген алгоритмдер. Башкача айтканда, машина татаал жана көп параметрдүү маселелерде (адам чече албаган) үлгүлөрдү таба алат, ошону менен так жоопторго жакын болгон жыйынтыктарды божомолдойт. Машиналык үйрөтүүнүн негизги максаты маалыматтар менен жыйынтыктардын ортосундагы закон ченемдүүлүктү туура аныктап, туура божомолдорду берүү менен адам баласынын жашоосун жакшыртуу жана жеңилдетүү. Азыркы учурда, машиналык үйрөтүүнү бизнес ээлери, маркетингдер жана кызматкерлер өз иштеринде туура чечим кабыл алыуу үчүн колдонууда. Киргизилген маалыматтарга негизделген эң так болжолдоолорду айтып берүү үчүн ар кандай алгоритмдер иштелип чыгууда. Машиналык үйрөтүүнүн натыйжасында машина жыйынтыкты алдын ала айтып, эстеп, керек болсо кайталап, бир нече варианттын ичинен эң жакшысын тандап алат.

Учурда машиналык үйрөтүү банктарда, ресторандарда, май куюучу станцияларда, өндүрүштөгү роботторго чейин кеңири чөйрөнү камтууда. Дээрлик күн сайын пайда болгон жаңы маселелер машиналык үйрөтүүнүн жаңы багыттарынын пайда болушуна алып келүүдө.

2.3. Машиналык үйрөтүүнүн түрлөрү

Машиналык үйрөтүүнүн методдору - гипотезаларды тастыктоого, жасалма интеллектти колдонуп оптималдуу чечимдерди чыгарууга багытталган маселелердин жыйындысы. Мугалимдин жардамы менен машиналык үйрөтүү(Supervised learning), Мугалимсиз машиналык үйрөтүү(Unsupervised learning), Жарым-жартылай мугалимдин жардамы менен машыиналык үйрөтүү(semi-supervised learning), Трансдуктивдик үйрөтүү (transductive learning), Үйрөнүүнү үйрөнүү (Learning to learn). Төмөндө машиналык үйрөтүү методдорунун ар бир түрүнүн реалдуу көйгөйлөргө карата колдонулушу тууралуу кененирээк маалымат берилет.

Мугалимдин жардамы менен машиналык үйрөтүү(Supervised learning) – машиналык үйрөтүүнүн бул түрү үйрөнүү учурунда мугалимдин жардамын талап кылат. Башкача айтканда, машиналык үйрөтүү үчүн даярдалган берилиштердин массиви адегенде туура жыйынтыктарды да камтыйт, андыктан алгоритмдин максаты туура жыйынтык берүү эмес, берилген берилиштердин жана туура жыйынтыктардын ортосундагы закон ченемдүүлүктү табуу аркылуу чечим кабыл алуу болуп саналат. Натыйжада туура болжолдоолорду жана моделдерди түзө алат. Бул машиналык үйрөтүүнү, тамга тааныбаган балага, мугалимдин тамга таанытуусу сыяктуу кароого болот.

Бул түр алгоритмди Фейсбук компаниясы сүрөттөрдү социалдык тармактарга жүктөө учурунда, сүрөттөгү башка адамдарды таанып, аларды белгилеп коюуну сушунш кылган учурда колдонгонун байкоого болот. Алгоритм ар бир киргизилген сүрөттү анализдеп ар бир адамды таанууну үйрөнөт жана үйрөтүлгөн алгоритм туура жыйынтыкттарды бере баштайт. Бул түр алгоритмди колдонуу социалдык тармактарда гана эмес медицина жаатында дагы кененирээк

колдонулганын билебиз. Мисалы: медицина тармагында онкологиялык оорулардын диагнозун коюуга ушул түр алгоритмдер колдонула баштаган. Берилиш катары оорулуу адамдардын ар түрдүү анализдеринин жыйынтыгы жана туура жыйынтык катары докорлордун койгон диагноздорун аркылуу үйрөтүү жүргүзүлгөн. Ушул сыяктуу берилиштин жыйынтыгын болжолдоо боюнча маселерди чечүүдө мугалимдин жардамы менен машиналык үйрөтүүнү колдонууга болот. Бул диссертациялык иштин максаты мугалимдин жардамы менен машиналык үйрөтүүгө багытталган. Тактап айтканда кошумча программанын чыгарып берген жыйынтыктары машиналык үйрөтүүгө кириш катары алынат жана мугалим тарабынан бааланган дил баяндын баалары жыйынтык болуп алынат. Мугалим дилбаянды баалоодо ар түрдүү критерийлерди колдонот, бул алгоритм дагы дал ошол критерийлердин жардамы менен ишке ашырылат. Алгоритмди үйрөтүү учурунда ар түрдүү темадагы жана ар кандай көлөмдө жазылган дил баяндар жана алардын мугалим тарабынан бааланган жыйынтыктары колдонулду. Алгоритмдин жалпы иштөө принцибинин схемасы Сүрөт 1де көрсөтүлгөн [9].

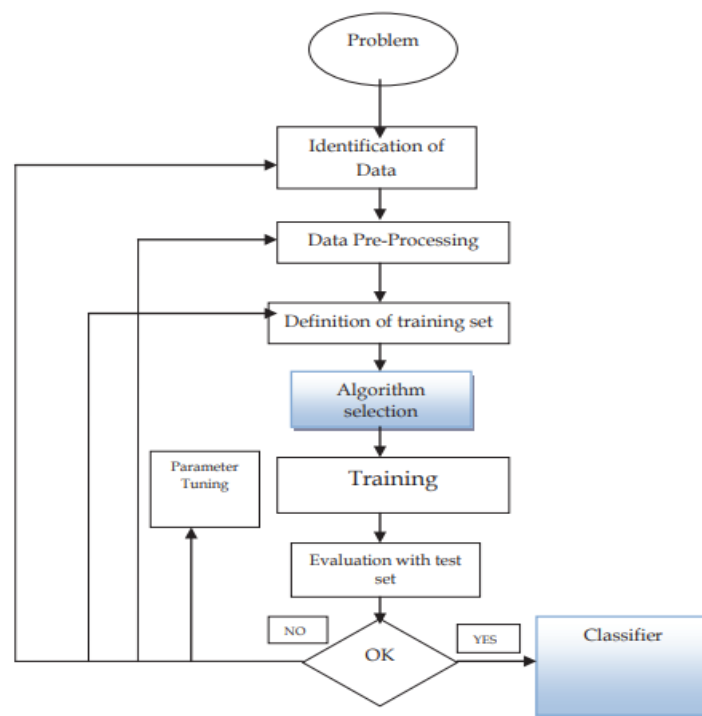


Fig. 2. Machine Learning Supervise Process

Сүрөт 1. Мугалимдин жардамы менен машиналык үйрөтүү схемасы.

Мугалимсиз машиналык үйрөтүү(Unsupervised learning). кутуунун бул түрү үчүн негизги түшүнүк үлгү болуп саналат - олуттуу маалымат массивдерин иштетүүдө алгоритм алгач өз алдынча үлгүлөрдү аныктоосу керек. Кийинки этапта, аныкталган үлгүлөрдүн негизинде, машина маалыматтарды чечмелеп, системалаштырат. Бул алгоритмдин түрү жыйынтык кандай боло турганын так айта албайт бирок, үлгүлөрдүн негизинде божомолдоп берет [10]. Бул түр алгоритмдер социалдык тармактарда көп колдонулат. Колдонуучулардын көргөн, уккан, окуган маалыматтарын анализдеп, кийин кандай маалымат алгысы келээрин божомолдоп, сунуштарын чыгарат. Мисалы YouTube видеохостингин карай кетсек, ал платформада көрүлгөн видеолор google аккаунтка байланышат жана видеолордун жаныры, убактысы, тибине карай сунуштарды бере баштайт. Кээ бир учурларда бир каналдын колдонуучуларынын көп бөлүгү ушул каналга катталган

деген маалымат берип, каналдарды да сунуш кыла баштады. Бул сыяктуу божомолдоолор көп учурда туура сунуштарды жасайт, бирок ошол эле учурда бир google аккаунт менен бир нече киши колдонсо, алгоритм аралаш сунуштарды бере баштайт, ошол учурда колдонуучулар эмне себептен андай болуп жатканын билбей таң калышы мүмкүн(жашоодон алынган мисал).

Терең үйрөтүү (Deep learning). Deep Learning чоң көлөмдөгү маалыматтар боюнча татаал эсептөөлөрдү жүргүзүү үчүн жасалма нейрон тармактарын көп катмардуулугун колдонот. Бул адамдын мээсинин түзүлүшүнө жана функциясына негизделген машина үйрөнүүнүн бир түрү. Терең үйрөтүү алгоритмдери мисалдардан үйрөнүү менен машиналарды үйрөтөт. Саламаттык сактоо, электрондук коммерция, көңүл ачуу жана жарнама сыяктуу тармактар көбүнчө терең үйрөнүүнү колдонушат. Азыркы күндө терең үйрөтүү алгоритмдери сүрөттөн же видео көрүнүштөрдөн объектерди таануу жана ажыратуу сыяктуу маселелерди чечүүдө колдонулууда. Машиналык үйрөтүүдө берилген маалыматтын жыйынтыгын билүү менен алгоритмдер түзүлсө, терең үйрөтүүдө нейрондук тармактын катмарлары колдонулуп, биринчи катмарда белгилүү гана маалыматты аныктаса, кийинки катмарларда андан да терең байкап, божомолдоолорду жүргүзөт. Мисалы сүрөттө кыз менен эркек баланын көрүнүшү тартылган болсо, терең үйрөтүүдө объекттердин өңүн, кырларын салыштырып айрымасын салыштыруу менен божомолдойт [11].

2.4. Машиналык үйрөтүүнүн алгоритмдери

Машиналык үйрөтүү киргизилген маалыматтарды кабыл алып, талдоочу программаланган алгоритмдерди колдонуп, андан кийин жарактуу диапазондон чыгуу баалуулуктарын божомолдойт. Жаңы маалыматтар киргизилгенде, бул

алгоритмдер алардын иштешин үйрөнүп, оптималдаштырып, убакыттын өтүшү менен өндүрүмдүүлүгүн жогорулатып, "интеллекти" өнүктүрүшөт.

Регрессия – мугалимдин жардамы менен машиналык үйрөтүүдө колдонулган ыкма болуп саналат. Бул үзгүлтүксүз өзгөрмө моделдөө жана болжолдоо үчүн колдонулат. Сызыктуу регрессия алгоритмин колдонуунун мисалдары болуп төмөнкүлөр саналат: кыймылсыз мүлктүн баасын болжолдоо, сатууну болжолдоо, студенттик экзамендердин жыйынтыгын болжолдоо, биржадагы акциялардын баасынын кыймылын болжолдоо. Регрессияда бизде берилиштер топтомун белгилейбиз жана чыга турган өзгөрмөнүн мааниси киргизилген өзгөрмөнүн маанилери менен аныкталат - демек, бул мугалимдин жардамы менен машиналык үйрөтүү ыкмасы. Регрессиянын эң жөнөкөй түрү сызыктуу регрессия болуп саналат, мында түз сызыкты берилиштер жыйындысына туура келтирүү аракети жасалат жана бул маалымат топтомунун өзгөрмөлөрүнүн ортосундагы байланыш сызыктуу болгондо гана ишке ашурууга мүмкүн.

Классификация - жасалма интеллектте жана машиналык үйрөтүүдө объекттердин топтомун формалдуу сүрөттөлүшүн талдоонун негизинде класстар деп аталган топторго бөлүү маселеси коюлган. Классификацияда байкоонун ар бир бирдиги кандайдыр бир сапаттык касиеттин негизинде белгилүү бир топко же класстарга блөлүштүрүлөт.

Классификация маселеси көптөгөн тармактарда колдонулат:

соодада - кардарларды жана товарларды классификациялоо маркетинг стратегияларын оптималдаштырууга жана чыгымдарды кыскартууга мүмкүндүк берет, телекоммуникация тармагында - абоненттердин классификациясы

лоялдуулуктун деңгээлин аныктоого, лоялдуулук программаларын иштеп чыгууга мүмкүндүк берет, медицинада жана саламаттыкты сактоодо - ооруларды диагностикалоо, рискке кирген топтор боюнча калкты классификациялоо, банк секторунда - кредиттик скоринг.

Кластерлөө - бул маалымат чекиттери боюнча топторду камтыган күчтүү машина үйрөтүү ыкмасы. Белгилүү маалымат чекиттеринин топтомун эске алуу менен, окумуштуулар ар бир маалымат пунктун өзүнчө топко классификациялоо же классификациялоо үчүн кластерлөө алгоритмин колдоно алышат. Теориялык жактан караганда, бир топко кирген маалымат чекиттери окшош мүнөздөмөлөргө же касиеттерге ээ. Башка жагынан алганда, өзүнчө топторго таандык маалымат чекиттери абдан уникалдуу мүнөздөмөлөргө же касиеттерге ээ болушат.

Өлчөмдүн азайышы. Андан ары визуалдаштыруу же жумушта колдонуу үчүн объект мүнөздөмөлөрүнүн массивин аз сандагы белгилерге кысуу. Мисалы, тармактар боюнча берилиш үчүн бир катар маалыматтарды архивдерге кысуу.

Аномалияларды издөө. Негизги массадан кыйла айырмаланган сейрек кездешүүчү жана адаттан тыш объектилерди издөө, мисалы, алдамчылык операцияларын издөө [10].

БӨЛҮМ III.

АДАБИЯТ ТАЛДОО

3.1. Баяндаманы текшерүүнүү автоматташтыруу тарыхы

Табигый тилде сүйлөө жана ой жүгүртүүсүн жазуу түрүндө билдирүү жөндөмү адамзаттын уникалдуу каражаты катары көрө алабыз, бирок табигый тилде формулалык же алгоритмдик спецификацияны табуу абдан оор маселе. Андыктан, жазууну баалаган компьютердик программаларды иштеп чыгуу аракеттери көп учурда ишке ашырылышы мүмкүн болбогон маселе катары каралышы таң калыштуу эмес [12]. Ошого карабастан, тил үйрөнүү деңгээлинде жазууну (текстти) баалоону автоматташтырууну, жазуу жөндөмдүүлүктөрүнүн негизин түзгөн лингвистикалык, морфологиялык жана сөз байлык сыяктуу жакшы жана жаман жазууну мүнөздөгөн көптөгөн өзгөчөлүктөрдүн негизинде ишке ашыруу мүмкүн деген ойлор көптөгөн изилдөөчүлөр тарабынан айтылып келмекчи.

Алгачкы изилдөөлөр, 1966 жылы Page аттуу изилдөөчү тарабынан сунушталган [13]. Бирок ал учурдагы технологиялыр, бул түр долбоорлорду аткарууга мүмкүнчүлүк берген эмес. 2003 жылы, табигый тилди иштетүү багытынын өнүгүүсү менен бирге студенттердин эссе тапшырмаларынын текшерүүсүн автоматташтыруу долбоорлор кайрадан изилдене баштаган. Мисалы, бул багытты сунуштаган Page жана изилдөөчү тобу, PEG (Project Essay Grading) долбоорун иштеп чыккан [14]. Долбоордо, жазууну баалоо үчүн сөздүн орточо узундугу, үтүрлүү чекиттердин саны жана сөз сейректиги жана башка ушул сыяктуу жазуу стилинин өзгөчөлүктөрү колдонулган. Бул өзгөчөлүктөр көптөгөн мугалим тарабынан текшерилген тапшырмалардан чыгартылып, кийин бул маалыматтын негизинде регрессия ыкмасын колдонуп, баалар божомолдонгон. Ал эми Intelligent Essay Assessor долбоорунун изилдөөчүлөрү, мурун топтолгон жазуу тапшырмалар

менен салыштыруу индекси аркылуу баалоо жүргүзүүнү сунуштаган. Дил баяндарды баалоону автоматташтыруу компьютер технологиялары өнүгө башатагдан бери кызыгуу жаратып келген. Дил баянды баалоо көп ресурстарды талап кылгандыгына байланыштуу, баалоону автоматташтыруунун үстүнөн көптөгөн илимий изилдөөлөр аткарылган. Ишке ашырылган илимий иштерди изилдөө учурунда коюлган максат такталып, теориялык билим толукталды. Ишке ашырылган моделдер: Project Essay Grade (PEG), Intelligent Essay Assessor (IEA), Educational Testing service I, Electronic Essay Rater (E-Rater), C-Rater, BETSY, Intelligent Essay Marking System, SEAR. Бул моделдер ар кандай максатта ишке ашырылган. Бул моделдердин арасында илимий изилдөө үчүн жана коммерциялык максатта ишке ашырылгандары дагы бар. Ар бир моделдин артыкчылыгы жана кемчилигин табууга болот. Салыштыруунун жыйынтыгында Educational Testing Service I модели 93-96 % га чейин туура баалаган [15].

3.2. Баяндаманы текшерүүнү автоматташтырууда колдонулган машина үйрөтүү алгоритмдери

Дил баяндарды баалоону автоматташтыруу учурунда, мугалим баалаган жана нейрондук тармак баалаган дилбаяндардын бааларын салыштыруу аркылуу баалоонун сапатын жакшыртууга болот. Баалоону салыштыруу үчүн ар кандай моделдер ишке ашырылган. Пирсондун корреляциясы, Спирмендин рангдык корреляциясы, Тау Кендаллдын Каппа Квадраттык Өлчөөсү (QWK- quadratic-weighted kappa). Мугалим жана нейрондук тармактардын баалоолорун салыштыруу аркылуу, машинанын канчалык деңгээлде туура баалоо жүргүзө ала тургандыгын билүүгө болот. Бул моделдердин жардамы менен машиналык окутуунун сапатын жакшыртууга жана автоматтык баа берүүнү өнүктүрүүгө болот [16].

Дил баяндарды баалоону автоматташтырууну эки негизги багытка бөлүүгө болот. Биринчи багыттагы моделдерди ишке ашырууда тилидин өзгөчөлүгү жана анын грамматикасынын туура колдонулушу каралат. Орфографикалык каталар, колдонулган сөздөрдүн саны, сүйлөмдөрдүн саны, грамматикалык каталар ж. б. критерийлер аркылуу түзүлүүчү моделдер. Экинчи багыттагы моделдерде эссенин маанисине көңүл бөлүнөт. Ушул сыяктуу моделдерди ишке ашырууда deep learning, терең үйрөтүү колдонулат [11].

Дил баянды баалоону автоматташтырууну ишке ашыруу үчүн машиналык үйрөтүү колдонулат. Машиналык үйрөтүүнүн бир нече түрлөрү бар жана алардын өзгөчөлүктөрүн салыштыруу аркылуу коюлган максатка жетүү үчүн ыңгайлуу болгону тандалып алынды.

Ишке ашырылган илимий ишиерди изилдөө учурунда, ар кандай жол менен түзүлгөн моделдер каралды.

Taghipour жана Ng илимий изилдөөсүндө дил баянды баалоону автоматташтыруунун моделин нейрондук тармакты колдонуу аркылуу ишке ашырышкан. Алгоритмди түзүү үчүн $s(x) = \text{sigmoid}(w \cdot x + b)$, сигмоиддик функцияны колдонушкан. Ал эми түзүлгөн моделди үйрөтүү үчүн ачык булактардан берилиштерди колдонушкан. Бул статьяда моделди үйрөтүү үчүн дилбаяндарды топторго бөлүп, жалпысынан 13000 ге жакын дил баян колдонулган. Жыйынтыгында $k = 0.805$ тактыкка чейин туура баалоого жетишкен [16].

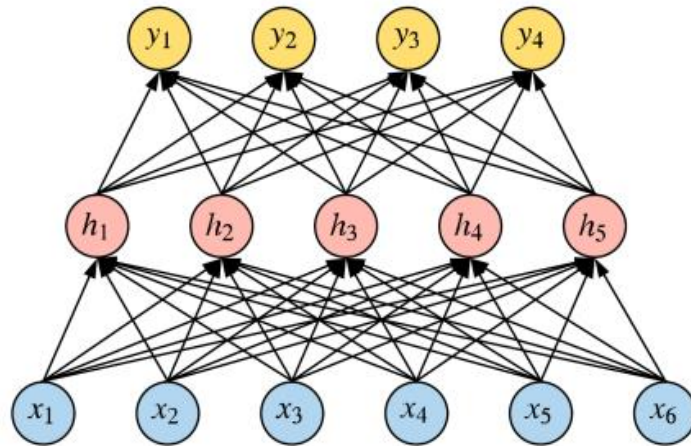
Song & Zhao изилдөөсүндө дил баянды баалоону автоматташтырууну машиналык үйрөтүү аркылуу иш жүзүнө ашырышкан. Бул изилдөөнүн максатын ишке ашырууда сызыктуу регрессия колдонулган.

Колдонулган сызыктуу функция:

$$h\theta(x) = \theta^T x + c$$

Ишке ашырылган моделдин эң жакшы жыйынтыгы $\kappa = 0.52$ ни көрсөткөн [17].

Дил баянды балоону автоматташтыруу долбоорлору көпчүлүк учурда англис тили үчүн ишке ашырылган жана англис тилинде жазылган эсселердин жыйнагын ачык булатардан алуу аркылуу моделди үйрөтө алабыз, үйрөтүлгөн моделди текшерүүгө да болот. Дил баяндарды автоматтык түрдөө баалоо көйгөйү Шведциянын улуттук сынактарында да пайда болгон. Улуттук сынактарда студенттер дил баян жазуу аркылуу баа топтогондугуна байланыштуу, бир канча миндеген студенттин дил баянын мугалимдер тарабынан текшерилүүсү абдан көп ресурсту талап кыла баштаган. Ушул көйгөйдү чечүү максатында дил баянды баалоону автоматташтыруу системасын иштеп чыгышкан. Иштелип чыккан системада мугалимдин жардамы менен машиналык үйрөтүү методу тандалган жана алгоритмди ишке ашырууда нейродук тармактарды колдонушкан. Нейродук тармактар төрт түрдүү жыйынтык бере тургандай кылып моделдешкен. Ар бир текст үчүн төрт мүмкүн болгон баа бар: IG (өтпөй калды), G (өттү), VG(артыкчылык менен өттү) жана MVG (эң жакшы). Сүрөт 2 де бул илимий изилдөөнүн нейродук караты көрсөтүлгөн. Бул эмгекте төрт түрдүү жыйынтык берээрин схема түрүндө көрүүгө болот [18].



Сүрөт 2. IG (өтпөй калды), G (өттү), VG (артыкчылык менен өттү) жана MVG (эң жакшы)

Бирок, бул багыттагы алгачкы изилдөөлөр 1966 жылда башталганга карабастан, көйгөй азыркы учурга чейин толугу менен чечилген эмес. Азыркы күндө бул темада көптөгөн илимий изилдөөлөр жүргүзүлүп жатат. Англис тили үчүн жасалган илимий изилдөөлөр жана алгоритмдер көп болгону менен түрк тили жана кыргыз тили үчүн жасалган алгоритмдер дээрлик жокко эсе. Бул изилдөөдө түрк тилинде жазылган дилбаяндарды баалоону автоматташтыруу ишке ашырылды. Бул алгоритмди кыргыз тилинде жазылган дилбаяндар үчүн да колдонууга болот.

БӨЛҮМ IV.

МЕТОД ЖАНА МАТЕРИАЛДАР

4.1. Метод жана материалдар

4.1.1. Тил китепканалары. Zemberek китепканасы

Түрк тилдүү элдер Европадан тарта Сибирге чейинки жерлерде жайгашкан жана 140 млн га жакын адам бул тилди колдонот. Түрк тил түркүмү тилдик өзгөчөлүктөрү боюнча окшош тилдердин тобу болуп саналат. Бул тилдерди өзгөчөлүгүн изилдөө максатында көптөгөн илимий иштер жарыяланган. Бирок көбүнчө эмгектер түрк тилин изилдөөгө багытталган. Ошол себептерден улам түрк тилинин тилдик өзгөчөлүгүн аныктаганга мүмкүндүк берген алгоритмдердин жыйнагы болгон Zemberek китепканасын ишке ашырылган. Китепканада түрк тилинин грамматикасын автоматтык түрдө текшерүү мүмкүнчүлүгү камтылган. Түрк тилинде жазылган жазууларды изилдөө максатында колдонууга ыңгайлуу. Ар бир тилдин тилдик өзгөчөлүгүн, ал тилде жазылган жазууларды изилдөө аркылуу табууга болот. Көптөгөн жазууларды текшерүү менен тилдин өзгөчөлүгүн камтыган статистикаларды чогултууга болот. Бул изилдөөлөр азыркы учурда өтө актуалдуу темалардын бири болууда.

Технологиянын өнүгүүсү менен тилди таанып билүү, башкача айтканда жазылган тексттин же аудио үндөрдүн кайсы тилге таандык экенин компьютердин жардамы менен аныктоо мүмкүнчүлүгү пайда болду. Мисалы: Google компаниясынын программдык жабдыгы болуп эсептелген Google котормочу текстти которууду, ал текст кайсы тилге таандык экенин көрсөтүү жөндөмдүүлүгүнө ээ. Бул автоматтык түрдө тилдү таанып билүү технологиясы ушул сыяктуу башка платформаларда дагы колдонулат. Речь же үндү таануу маалымат издөө учурунда дагы көп

колдонулат. Акылдуу телефондордо жардамчы кызматтардын көбү үндү таануу алгоритмдерин колдонот. Ушул сыяктуу алгоритмдер машиналык үйрөтүү технологиясын колдонуп ишке ашырылат.

Zemberek китепканасы түрк тилинде жазылган дилбаяндарды баалоону автоматташтыруудагы биринчи этап болуп саналат. Китепкана түрк тилинин тилдик өзгөчөлүгүн камтып, дилбаяндын грамматикалык жана орфографиялык каталарын табууга жардам берет. Дилбаянда колдонулган ар бир сөздү сөз түркүмүнө ажыратып берет. Сүйлөмдүн санын, ката колдонулган сөздөрдүн санын, сөз байлыгын, этиштин ыңгайларынын, чактарынын, жактарынын санын көрсөтүп берет. Мугалим дилбаянды баалап жаткан учурда сүйлөмдө кездешкен ар бир сөздүн сөз түркүмүнө маани бербейт, бирок сөз түркүмдөрүнүн туура колдонулушу өтө маанилүү. Дилбаянда кездешкен сөз түркүмдөрүнүн саны дилбаянды жазган адамдын тилди канчалык деңгээлде өздөштүргөнүн билүүгө болот. Эгер окуучу этиш жана зат атоочтук сөздөр менен гана сүйлөм түзө алса, анда окуучу тилдин баштапкы этабын гана өздөштүрдү деп баа берүүгө болот. Башка сөз түркүмдөрүн колдонуп сүйлөм түзө баштаганда окуучунун сөз байлыгы өнүккөнүн байкоого болот. Ушул сыяктуу божомолдорду колдонуу аркылуу дилбаянды баалоону автоматташтырууга болот.

Zemberek китепканасы Java программалоо тилинде ишке ашырылган. Дилбаянды баалоону автоматташтырууда колдонулган программдык жабдык Java программалоо тилинде проектилеген. Java программалоо тили кеңири таралган тилдердин бири болгондугуна байланыштуу, кодоо учурунда пайда болгон көйгөйлөрдү чечүү жеңил болду. Программаны колдонуу ыңгайлуу болсун үчүн интерфейс түзүлгөн. Интерфейске бир нече дилбаянды бир учурда жүктөөгө болот

жана бир файлга бардык жыйынтыктарды жүктөп алуу мүмкүнчүлүгүнө ээ. Жүктөлгөн жыйынтыктар андан кийинки машиналык үйрөтүүдө колдонулат.

4.1.2. Колдонуучу интерфейсти иштеп чыгуу

Ишке ашырылган программдык жабдыктын интерфейсін Сүрөт 3 тө көрүүгө болот. Программдык жабдыкка бир эле учурда бирдей темадагы бир нече дилбаянды киргизип, текшертип алуу мүмкүнчүлүгү бар. Берилген бардык дилбаяндын жыйынтыгын excel файлга топтойт. Андан кийинки кадамда керек болгон маалыматтар машиналык үйрөтүү үчүн колдонулат.

Жардамчы программдык жабдык чыгарган жыйынтык төмөндө көрсөтүлдү:

Kullandığı Anahtar Kelimler : komedi gülmek sorun

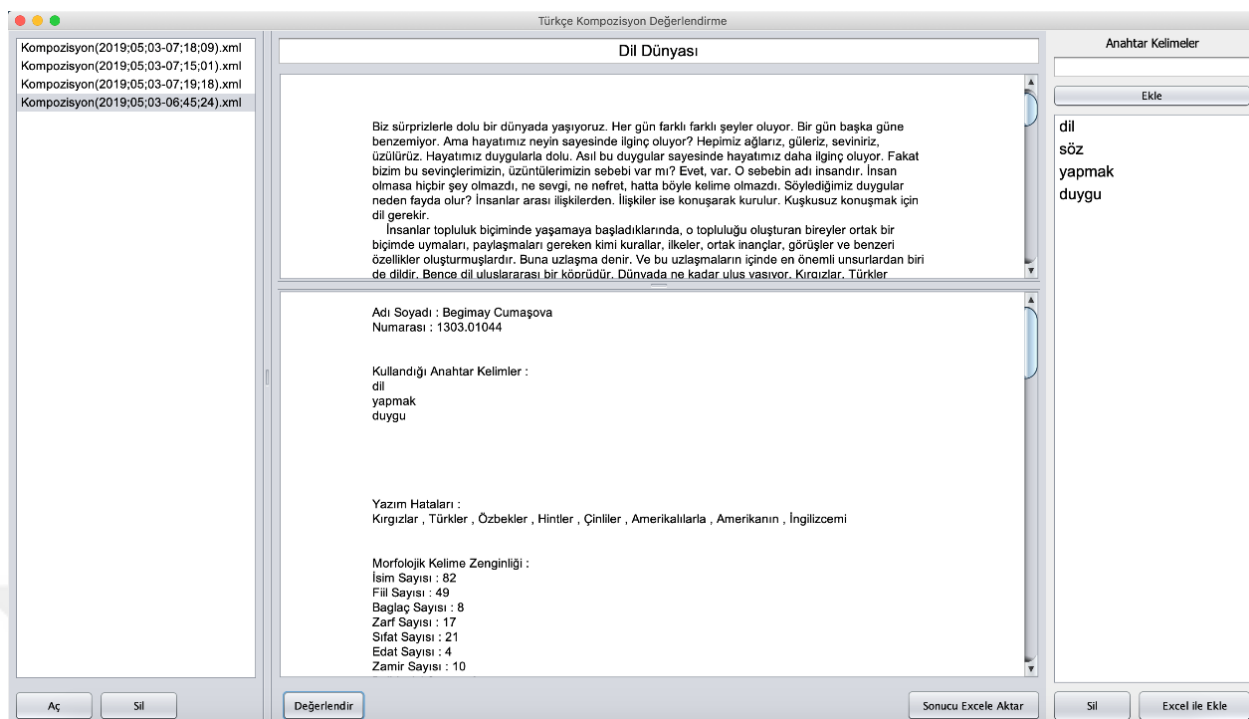
Yazım Hataları :39

Bütün , dı , hokabaz , sehatliğe , Herkez , hokabazı , istiyorlardı , istemediğini , hokabazlıkla , değişverdi , Orhana , Hokabaz , hokabazımız , başarısızlığa , vardu , herhengi , binzemiyor Morfolojik Kelime Zenginliği : İsim Sayısı : 79 Fiil Sayısı : 63 Bağlaç Sayısı : 7 Zarf Sayısı : 24 Sıfat Sayısı : 24 Edat Sayısı : 3 Zamir Sayısı : 11 Belirleyici Sayısı : 5 Noktalama Sayısı : 5 Özel İsim Sayısı : 7 Sayıların Sayısı : 3 Bilinmeyenlerin Sayısı : 1 Kelime Zenginliği : İsim (79) bugün , sınıf , film , ad , hokkabaz , komedi , yan , ora , olay , fil , baş , çocuk , arkadaş , iskender , iş , yardımcı , göz , ameliyat , para , turne , karar , enişte , karavan , baba , şehit , dâhil , yolcu , yol , bey

, şey , komik , köy , düşün , program , sahne , gelin , akşam , oyun ,
sıra , seyirci , bura , sihir , kutu , iç , kimse , adam , akıl , soru , aşk ,
kız , sevgili , insan , anne , ağabey , oğul , mektup , devam , vb , gün
, koltuk , hayat , zevk , su , ara , ev , nefes , sürpriz , bayan , sonuç ,
stres , numara , gerek , var , zaman , dakika , yardım , dünya , sorun
, hediye *Fiil (63)* izlemek , gülmek , düşünmek , anlamak , görmek ,
büyümek , olmak , yapmak , istemek , kazanmak , çıkmak , vermek ,
gitmek , yaşamak , etmek , kalmak , sahnelemek , gelmek , hazırlamak
, tanışmak , başlamak , göstermek , girmek , kaybolmak , şaşırmak ,
suçlamak , demek , kaçmak , evlenmek , zorlamak , değil ,
karşılaşmak , bilmek , ölmek , beğenmek , almak , söylemek ,
değişmek , kanmak , bırakmak , sevmek , uğraşmak , bakmak ,
yazmak , oturmak , geçirmek , kaybetmek , düzenlemek , kalkmak ,
unutmak , satmak , koşmak , açmak , üzme , tutmak , sevinmek ,
bitmek , benzetmek , uğramak , güvenmek , desteklemek , benzemek
, getirmek *Sıfat (24)* geçen , beraber , zor , ilginç , iyi , büyük , lazım ,
ne , böyle , var , basit , doğru , son , dolu , yeni , kapanık , ilk ,
hokkabaz , yakın , herhangi , kötü , yalnız , gerçek , zengin *Edat (3)*
için , kadar , diye *Zarf (24)* çok , sonunda , sonra , ne , evet , hemen
, şöyle , acaba , belki , aslında , sadece , hep , herhâlde , aniden ,
artık , ise , yeniden , hiç , yine , böylece , kısacası , daha , ancak , işte

Zamir (11) biz , o , biri , herkes , ben , kim , kendi , bu , ne , sen , şu
Özel İsim (7) orhan , çanakale , iskender , sait , çanakkale , fatma ,
başar Sayı (3) iki , on , ikinci Belirlenemeyen (1) d Haber
Kipleri(Zamanlar) : GEÇMİŞ ZAMAN : BELİRSİZ GEÇMİŞ ZAMAN:
Belirsiz Geçmiş Zaman : 3 Belirsiz Geçmiş Zamanın Hikayesi : 0
Belirsiz Geçmiş Zamanın Rivayeti : 0 Belirsiz Geçmiş Zamanın Şartı
: 0 BELİRLİ GEÇMİŞ ZAMAN: Belirli Geçmiş Zaman : 13 Belirli
Geçmiş Zamanın Hikayesi : 0 Belirli Geçmiş Zamanın Şartı : 0
ŞİMDİKİ ZAMAN: Şimdiki Zaman : 62 Şimdiki Zamanın Hikayesi : 3
Şimdiki Zamanın Rivayeti : 4 Şimdiki Zamanın Şartı : 0 GENİŞ
ZAMAN: Geniş Zaman : 7 Geniş Zamanın Hikayesi : 1 Geniş
Zamanın Rivayeti : 0 Geniş Zamanın Şartı : 0 GELECEK ZAMAN:
Gelecek Zaman : 4 Gelecek Zamanın Hikayesi : 0 Gelecek Zamanın
Rivayeti : 0 Gelecek Zamanın Şartı : 0 Dilek Kipleri : ŞART KİPİ: Şart
Kipi : 1 Şart Kipinin Hikayesi : 0 Şart Kipinin Rivayeti : 0 EMİR KİPİ:
Emir Kipi : 2

GEREKLİLİK KİPİ: Gereklilik Kipi : 0 Gereklilik Kipinin Hikayesi
: 0 Gereklilik Kipinin Rivayeti : 0 İSTEK KİPİ: İstek Kipi : 1 İstek Kipinin
Hikayesi : 0 İstek Kipinin Rivayeti : 0 Kullandığı Toplam Kelime Sayısı
: 465 Toplam Cümle Sayısı : 75



Сүрөт 3. Ишке ашырылган программдык жабдыктын интерфейси

4.1.3. Машина үйрөтүүдө коллонулган технологиялар

Python программалоо тили илим изилдөөдө, илимий эсептөөлөрдө популярдуу тилдердин бири болуп калды. Жогорку деңгээлдеги интерфейси жана илимий китепканалардын көбөйүүсү бул тилде алгоритмдерди иштеп чыгуу үчүн алштырылгыс инструменттердин бири кылды. Тил илимий изилдөөдө гана эмес, өндүрүштө да колдонула баштады. Бул программалоо тили статистикалык маалыматтарды талдоодо жана топтолгон маалыматтарды машиналык үйрөтүүдө кеңири колдонулууда. Бул тилди табигый илимдерди талдоодо гана колдонбостон, гуманитардык илимдерди да изилдөөдө колдонулууда. Бул диссертациялык иште табигый тилди талдоо жана анда жазылган дилбаяндарды баалоо үчүн Python программалоо тили жана анын китеп канасы болгон Numpy колдонулду. Numpy: моделдин маалыматтары жана параметрлери үчүн колдонулган негизги маалымат

структурасы. Киргизилген маалыматтар башка илимий Python китепканалары менен интеграцияланууну жеңилдеткен чексиз массивдер катары берилет. Numpy'дын көз карашка негизделген эс тутум модели компиляцияланган код менен байланышканда да көчүрмөлөрдүн санын чектейт. Ал ошондой эле негизги арифметикалык амалдарды да камсыз кылат. Numpy китепканасынан тышкары excel программасы менен байлашуу үчүн xlrd китепканасынын мүмкүнчүлүктөрү колдонулду. Xlrd китепканасы excel файлында сакаталган берилиштерди туура окуп алуу максатында колдонулду.

4.2. Изилдөө кадамдары

Изилдөө алкагында, табигый тилде жазылган тексттердин өзгөчөлүктөрүн чыгаруу алгоритмдерин колдонуп, жыйынтыктардын негизинде баалоо жүргүзүлдү. Бул багытта жасалган изилдөөлөрдө, автоматташтырылган дилбаянды баалоонун үч ыкмасын, атап айтканда Latent Semantic Analysis (LSA) [5-7] жана буга тиешелүү Probabilistic LSA (PLSA) статистикалык моделдерди [8-9], жана Latent Dirichlet Allocation (LDA) [10-11] түр ыкмалар колдонулат. Бирок бул түр ыкмалар азыркы учурга чейин англис тили үчүн колдонулуп келген. Ошондуктан диссертациялык иште, түрк тилине ыңгайлуу ыкма тандалат.

Баяндама тапшырмаларын баалоону божомолдоо үчүн өрнөк тапшырмалардын жыйыны колдонот. Изилдөөнүн чоң тобу (максат болгон масса) – бул тил үйрөнгөн студенттердин баян түрүндөгү тапшырмалар, ал эми кичине тобу (үлгү болгон масса) – божомолдоо моделин иштеп чыгуу үчүн колдонулуучу текшерилген баяндамалардын тобу болот. Жогоруда белгиленгендей, изилдөөнүн кичине тобу (үлгү болгон масса) – божомолдоо моделин иштеп чыгуу үчүн колдонулган текшерилген баяндамалар болот. Бул маалыматты, институттун уруксаты менен Университетибиздин Чет тилдер жогорку мектебинин Тил үйрөтүү бөлүмүнөн

алынды. Этикалык нормаларды сактоо үчүн баяндамалардын студенттердин аты-жөнү өчүрүлгөн версиялары алынды.



БӨЛҮМ V.

БЕРИЛИШТЕРДИ ДАЯРДОО ЖАНА МАШИНАЛЫК ҮЙРӨТҮҮ

5.1. Берилиштерди даярдоо

Берилиштер катары, Кыргыз-Түрк Манас университетинин тил окутуу бөлүмүндө түрк тилин үйрөнүп жаткан студенттердин сабакта жазган дилбаяндары колдонулду. Жалпысынан, 50 темада жазылган 10 000 дилбаян топтолгон. Алар машина үйрөтүү көптүгү жана тестрлөө көптүгүнө бөлүндү. Машина үйрөтүү көптүгү катары 5000 дилбаян колдонулуп, ал эми жыйынтыктарды тастыктоого бөлүнгөн тестрлөө көптүгү 5000 дилбаяндан түзүлдү.

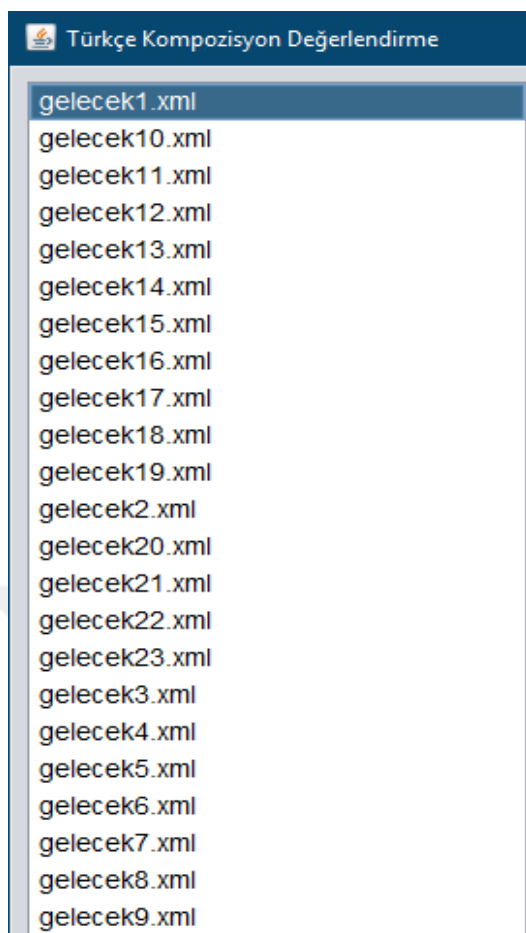
Берилиштер, этикалык эрежелер сакталыш үчүн анонимдүү түрдө алынып, жалгыз гана дилбаян өзү жана мугалим берген баасы колдонулду.

Машина үйрөтүү үчүн берилиштер атайы структурага келтирилди. Алынган дилбаяндар .xlsx форматына келтирилип топтолгон. Ал форматтагы маалыматты .xml форматына өткөрүү үчүн атайын кошумча программдык жабдыкка кошумча код кошулду. Табигый тил иштетүүчү кошумча программдык жабдык .xlsx форматындагы маалыматтарды .xml форматына өткөрүп сактайт. .xml форматына өткөрүлгөн дилбаяндын көрүнүшү Сүрөт 4 тө көрсөтүлгөн. .xml форматындагы ар бир тегдин өзүнүн мааниси бар, ТТИ кошумча программдык жабдыктын алгоритми тегдерди ачык катары карайт жана тег камтыган маалыматтарды ачык тегдерге жараша бөлүп алып иштетет.

```
<?xml version='1.0' encoding='UTF-8'?>
<kok>
<baslik>Gelecek</baslik>
<icerik>Merhaba. Size geleceğimi veya nasıl gördüğümü anlatmak istiyorum. Gelecekte
üniversiteyi mükemmel bir şekilde bitirmek ve profesyonel olmak istiyorum. Sonra seyahat
etmek ve kendim ve kariyerim için yeni bir şey aramak istiyorum. Dünyanın her yerinden
dolaşmak istiyorum. İş ve kariyerimi kurmak için Amerika'ya ve Amerika'ya taşınmak
istiyorum. Gezegendeki en iyi biyoteknolog olmayı hayal ediyorum. 25 ya da 26 yaşında
evleneceğim. Ailemin, karımın, çocuklarımın, kız kardeşlerimin parlak bir geleceği ve
uzun, mutlu bir hayatı olmasını istiyorum. Sanırım bir sürü çocuğum olacak ve bunu
istiyorum. Düşünüyorum ve her zaman ailemle birlikte olmak istiyorum. Benim için
hayattaki en önemli şey sevdiklerimin sağlığı ve sonra ikinci sırada kariyer ve iş. İçimde
değişecek çok şey olduğunu düşünüyorum, örneğin karakterim ve iç dünyam. Her insanın
büyüme döneminde farklı, birisi daha iyi bir şekilde değişecek, diğerleri daha iyi değil,
ama daha iyi bir şekilde değişmeye ve iyi bir insan olmaya çalışacağım. Büyük yazar
Cengiz Aytmatov'un dediği gibi, " Ama her insanın önünde, bugün, yarın, her zaman insan
olmak için sonsuz bir görev var.", benim için de en önemli şey her zaman bir insan olmak
ve bir hayvan değil. Bununla herkese hikayeler anlatmak istiyorum, sadece her zaman insan
olarak kalın.</icerik>
</kok>
```

Sürət 4. xml formatına ötkөрүлгөн дилбаяндын көрүнүшү

.xml formatına ötkөрүлгөн бардык дилбаяндар бир учурда кошумча программдык жабдыктын интерфейси аркылуу киргизилип, чыккан жыйынтыктарды .xlsx formatына сактайт. Сүрөт 5те .xml formatына ötkөрүлгөн дилбаяндардын массиви кошумча програмдык жабдыктын интерфейсінде көрсөтүлгөн. Бир учурда киргизилген дилбаяндардын массиви бир текшерилип, бир .xlsx formatындагы файлга жүктөлөт.



Сүрөт 5. .xml форматына өткөрүлгөн дилбаяндардын массиви

Сүрөт 6 да .xlsx форматында саталган жыйынтыктардын көптүгү көрсөтүлгөн. Ушул жыйынтыктар машиналык үйрөтүү үчүн курулган нейрондук тармакты моделдөөдө кириш маалымат же берилиш катары колдонулду.

Yazım Hataları sayısı	Isimler	Fiiller	Baglaclar	Zarflar	Sifatlar	Edatlar	Noktalam	Ozelisimler	Sayılar
38	36	23	4	6	12	1	3	1	0
0	18	14	1	2	6	4	2	1	0
13	32	16	4	4	7	3	2	6	1
0	33	13	2	3	6	4	2	1	0
6	27	14	3	3	4	2	1	3	4
0	31	13	2	4	7	5	2	1	2
1	36	21	3	6	9	3	2	13	2
1	20	13	6	4	7	1	2	3	1
7	31	28	6	8	13	3	7	2	3
4	54	29	6	11	15	3	3	6	0
1	40	19	6	4	10	4	3	3	3
1	30	18	5	6	7	5	4	5	1
2	29	27	2	4	8	1	2	2	3
0	22	18	4	5	7	5	2	0	0
4	65	29	6	4	15	5	3	9	6
1	37	15	3	7	5	3	3	3	1
1	37	22	4	5	8	4	2	8	0
3	24	11	2	2	7	3	2	1	0
3	23	8	2	4	2	2	2	5	3
4	46	19	2	4	11	4	2	4	0
9	25	14	5	2	4	2	2	5	0
29	21	14	1	4	4	5	3	15	3

Сүрөт 6. .xlsx форматында саталган жыйынтыктар

5.2. Машиналык үйрөтүү процесси

Машиналык үйрөтүүнү ишке ашыруу үчүн нейрондук тармактарды колдонуу аркылуу модель түзүлдү. Түзүлгөн модел мугалимдин жардамы менен машиналык үйрөтүү түрүн колдонот. Буга чейинки ишке ашырылган системанын жыйынтыгы нейрондук тармактарга берилиш катары киргизилет жана үйрөнүү процесси ишке ашырылат. Нейрондук тармакты ишке ашырууда активация кылуучу функциялар колдонулат. Активация кылуу функциялары: сигмоида, түз сызыктуу, тепкичтүү, ReLu, tanh жана башкалар. Сигмоида функциясынын так чектери бар $[0,1]$, ал эми түз сызыктын чектери чексиз болгондуктан анын жыйынтыктары туруксуз болот. Нейрондук тармактар аркылуу божомолдоо жүргүзүүдө так чектүү активдештирүү функциясын колдонуу ыңгайлуу. Ошондуктан моделди ишке ашырууда Сигмоиддик функция колдонулду [19]. Колдонулган функциянын формуласы:

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

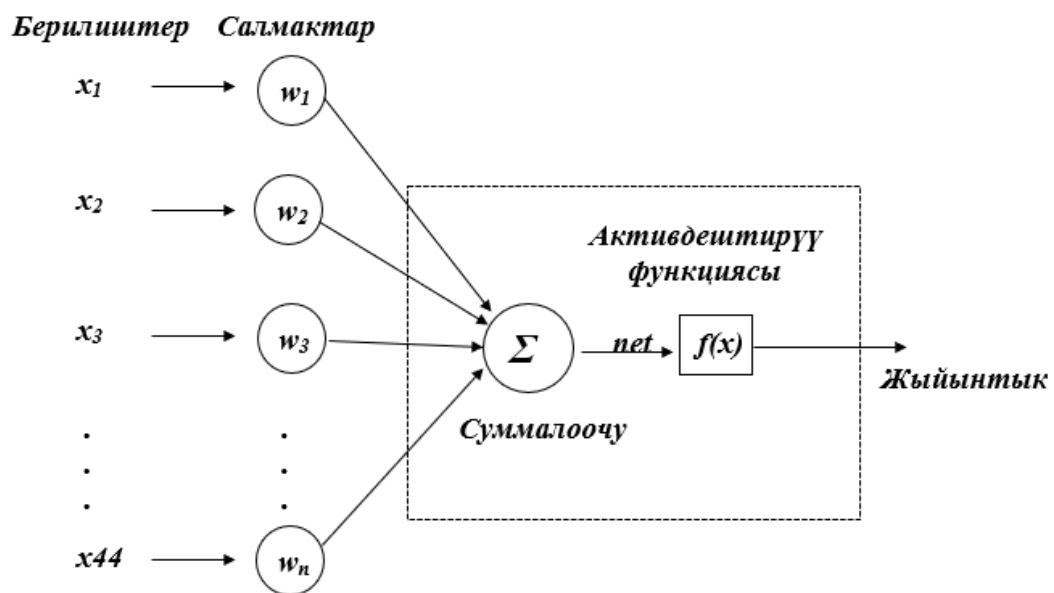
Бул формуланын берилиши катары x белгисизи көрсөтүлгөн. Активдештирүү функциясына:

$$hI = f_{activation}((x_1 * w_1) + (x_2 * w_2) + (x_3 * w_3) + \dots + (x_n * w_n))$$

берилиштер менен салмактардын арасында жүргүзүлгөн амалдардын жыйынтыгы жөнөтүлөт. x_n - бул нейрондук тармакты түзүү үчүн керек болгон кириш берилиштердин көптүгү, ал эми w_n - нейрондук тармактарды байланыштырган нейрондун салмактары. Нейрондун салмактары рандом түрүндө алынат. Ар бир үйрөнүү этабында нейрондун салмактары өзгөрүүгө учурайт.

Нейронду активдештирүүнү ишке ашыруу үчүн функциянын туундусун алуу менен ишке ашырылат

$$f'(x) = f(x) * (1 - f(x)),$$



Сүрөт 7. Нейрондук тармакты ишке ашыруу схемасы

Нейрондук тармактар нейрондун салмактарынан жана чектеринен түзүлөт. Биология жана медицина илиминде биологиялык нейрондорду бири-бири байланыштырган жипчелерди синапс деп аташат. Ошол эле терминдер жасалма интеллектти ишке ашурууда дагы колдонулат. Нейрондук тармактар табигый нейрондор сыяктуу моделденет. Башкача айтканда жасалма нейрондорду моделдөөдө табигый нейрондордун иштөө принциптери колдонулат. Синапс бул жасалма интеллекте нейрондорду байланыштырган, нейрондордун салмактары. Нейрондун салмактары берилиштерди башка нейрондорго өткөрүп жаткан учурда маалыматты алмаштырууга жөндөмдүү. Нейрондук тармактардын жардамы менен машиналык үйрөтүүнү моделдөөдө берилиштер менен жыйынтыктардын ортосундагы байланышты табуу үчүн салмактарды алмаштырып, оптималдаштырып, жоготууларды эсептеп, божомолдойт.

Нейрондук тармактын салмактары жана чектери кокус сандар аркылуу np.random.normal функциясын жардамы менен алынат. Ар бир этапта жоготуулар азаят. Баалоонун сапаты, салмактын жана чектердин жакшыруусунан түздөн түз көз каранды. Эпоха же этап деген түшүнүк окутуу учурунда канча жолу бардык берилиштер топтому аркылуу кайталап өтүүнү айтууга болот. Ар бир эпохада ар башка салмактар алынып, ар башка чектер берилет. Ошого жараша каталар дагы азаят. Сүрөт 8 де эпоханын өсүүсү менен жоготуулардын азаюусун байкоого болот.

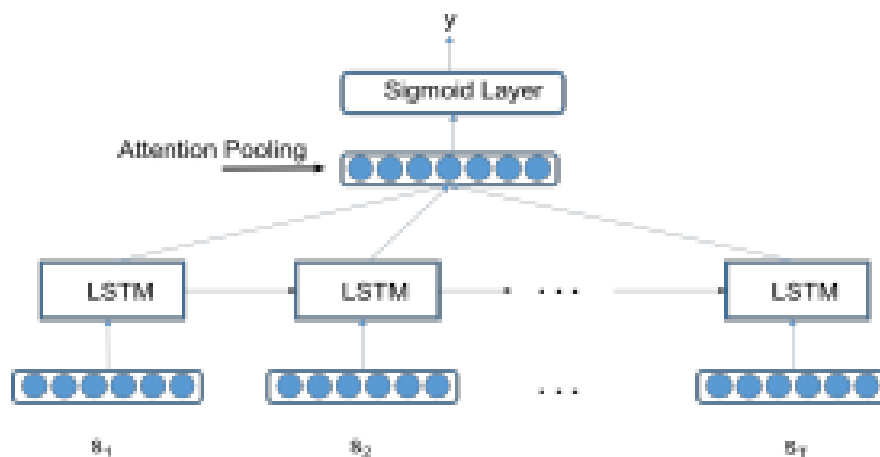
```

Epoch 0 loss: 0.058
Epoch 100 loss: 0.029
Epoch 200 loss: 0.027
Epoch 300 loss: 0.027
Epoch 400 loss: 0.027
Epoch 500 loss: 0.027
Epoch 600 loss: 0.027
Epoch 700 loss: 0.027
Epoch 800 loss: 0.027
Epoch 900 loss: 0.027
Epoch 1000 loss: 0.027

```

Сүрөт 8. Жоготуулардын эпохадан көз карандылыгы

Нейрондук тармактын модели Сүрөт 9 да схема түрүндө көрсөтүлгөн. Бул схемада бир нече кириш маалыматты жана бир чыгыш маалыматты камтыган нейрондук тармактардын моделин көрүүгө болот .



Сүрөт 9. Нейрондук тармактын модели

Сүрөт 9 да көрсөтүлгөн $S(n)$ нейрондук тармактын берилиштери (input). Ал эми LSTM нейрондун эс тутуму, Sigmoid Layer нейрондук тармакты активдештирүү

функция болуп саналат. Курулган модел бир гана жыйынтык берет жана ал жыйынтык 0 жана 1 санынын ортосундагы каалагандай бир санды берет.

Дилбаянды текшерүүгө керек болгон алгоритм толук иштелип чыкты жана алынган жыйынтыктар топтолду. Топтолгон жыйынтыктар машиналык үйрөтүүдө колдонулду. Дилбаянды автоматтык түрдө текшерүү системасынын схемасы Сүрөт 10 до көрсөтүлгөн. Схемадан кошумча программдык жабдык менен нейрондук тарматын байланышын жана машиналык үйрөтүү кандай жол менен ишке ашырылганын көрүүгө болот. Блок схемада “үйрөтүүчү берилиштер” блогу нейрондук тармакты машиналык окутуу үчүн берилген дилбаяндардын коптугу. Ал эми “Дилбаян берилиштери” блогу окутуп, үйрөтүлгөн нейронду текшерүү катарында берилген дилбаяндардын массиви. “Жакшыртуу үчүн сунуш”- бул кетирген каталарды эсептөө менен нейрондун салмактарын алмаштыруу керек экендигин көрсөткөн блок [20].

ТЕКСТ ТАПШЫРМАЛАРЫН АВТОМАТТЫК ТҮРДӨӨ БААЛОО СИСТЕМАСЫ



Сүрөт 10. Дилбаянды автоматтык түрдө текшерүү системасынын схемасы

Өрнөк баяндамалардын баалоо жыйынтыктарынын негизинде машина үйрөтүү методунун жардамы менен модель курулду. Корутундулоо кадамында, иштелип чыккан модель, програмдык жабдык же плагин түрүнө капсуляцияланат. Табигый тилди изилдөө үчүн алгоритмдерди иштеп чыгуу абдан оор көйгөйлөрдүн бири болуп саналат. Ошондуктан табигый тилде жазылган баяндаманы дагы талдоо жана баалоо иш-чаралары кандайдыр бир тактыкка чейин гана аткарылышы мүмкүн. Албетте, компьютердин жардамы менен баяндама өзгөчөлүктөрүн чыгаруу – бул баяндаманы туура баалоо үчүн жардамчы аспап боло алат, бирок табигый тил татаалдыгы, мугалим сыяктуу толук кандуу баа берүүгө мүмкүнчүлүк бербесин унутпашыбыз керек.

БӨЛҮМ VI.

ЖЫЙЫНТЫКТАР

6.1 Машина үйрөтүү процессинде алынган алгач жыйынтыктар

Сүрөт 5 те кошумча программдык жабдык дилбаяндарды текшергенден кийин чыгарып берген жыйынтыктардын массиви. Бул жыйынтыктар машиналык үйрөтүүгө кириш берилиш катары кирди жана машиналык үйрөтүү ишке ашырылды. Машиналык үйрөтүү үчүн ушул сыяктуу дилбаяндардан 10 000ге жакыны колдонулду. Колдонулган дилбаяндар темаларга жана коюлган бааларына, дилбаяндын өлчөмүнө жараша класстарга ажыратылып, машианлык окутуу жүргүзүлдү. Машиналык үйрөтүү үчүн нейрондук тармак моделденди жана кошумча программдык жабдык чыгарып берген жыйынтыктардын көптүгү нейрондук тармактар үчүн берилиш катары колдонулду. Кошумча программдык жабдык дилбаяндарды текшерип, жыйынтыгында ал дилбаянда кездешкен сөздөрдүн, сүйлөмдүн, ката жазылган сөздөрдүн, ар бир сөз түркүмүнө тиешелүү болгон сөздөрдүн, этиштин ыңгайынын, этиштин чагынын, этиштин жагынын санын чыгарып берет. Дилбаяндын синтаксистик жана морфологиялык талдоосунун жыйынтыгы 44 берилиштен турган массивди түздү. Топтолгон массив нейрондук тармак үчүн кирүү берилиштери болуп тандалды. Сүрөт 11 жана Сүрөт 12 де кошумча программдык жабдык чыгарып берген жыйынтыктардын массиви көрсөтүлгөн. Ар бир дилбаяндын мугалим баалаган баасы жана кошумча программдык жабдык аныктап берген берилиштердин (колдонулган сөз жана сүйлөм, сөз түркүмдөрү, гамматика жана орфография) аталышы жана мааниси жазылган.

Grade/100,00	Dogru kelimele	Yazım Hatalar	Isimler	Fiiller	Baglaci	Zarfilar	Sifatlar	Edatlar	Noktalar	Ozelisim	Sayilar	Bilinmey	Belirsiz	(Belirsiz	(Belirsiz	(Belirsiz	(Belirli	Gr	Belirli	Gr	Belirli	Gr	Simdiki	Simdiki	Simdiki	Simdiki
48,00	0,74	29	21	14	1	4	4	5	3	15	3	7	0	0	0	0	0	0	0	0	0	10	0	0	0	
78,00	1,49	1	40	19	6	4	10	4	3	3	3	2	0	0	0	0	0	0	0	0	0	12	0	0	0	
82,00	1,80	1	36	21	3	6	9	3	2	13	2	4	0	0	0	0	0	0	0	0	0	10	0	0	0	
80,00	1,72	7	31	28	6	8	13	3	7	2	3	4	0	0	0	0	2	0	0	0	14	0	0	0		
28,00	0,73	0	18	14	1	2	6	4	2	1	0	7	0	0	0	0	0	0	0	0	3	0	0	0		
52,00	1,00	6	27	14	3	3	4	2	1	3	4	1	3	0	0	0	0	0	0	0	9	0	0	0		
60,00	0,95	2	29	27	2	4	8	1	2	2	3	1	0	0	0	0	0	0	0	0	5	2	0	0		
64,00	1,36	13	32	16	4	4	7	3	2	6	1	5	0	0	0	0	0	0	0	0	10	0	0	0		
56,00	0,92	1	20	13	6	4	7	1	2	3	1	2	0	0	0	0	0	0	0	0	3	0	0	0		
78,00	1,78	4	46	19	2	4	11	4	2	4	0	3	0	0	0	0	0	0	0	0	9	0	0	0		
76,00	1,48	1	37	15	3	7	5	3	3	3	1	3	0	0	0	0	0	0	0	0	9	0	0	0		
62,00	1,07	3	24	11	2	2	7	3	2	1	0	2	0	0	0	0	0	0	0	0	2	0	0	0		
60,00	1,15	0	22	18	4	5	7	5	2	0	0	1	0	0	0	0	0	0	0	0	7	0	0	0		
40,00	1,30	9	25	14	5	2	4	2	2	5	0	1	0	0	0	0	0	0	0	0	3	0	0	0		
72,00	1,24	0	31	13	2	4	7	5	2	1	2	1	0	0	0	0	0	0	0	0	12	0	0	0		
40,00	1,16	3	23	8	2	4	2	2	2	5	3	2	0	0	0	0	0	0	0	0	8	0	0	0		
60,00	1,43	38	36	23	4	6	12	1	3	1	0	4	0	0	0	0	0	0	0	0	13	0	0	0		
60,00	1,14	0	33	13	2	3	6	4	2	1	0	1	0	0	0	0	0	0	0	0	1	0	0	0		
60,00	1,17	1	30	18	5	6	7	5	4	5	1	3	0	0	0	0	1	0	0	0	6	0	0	0		
90,00	2,55	4	54	29	6	11	15	3	3	6	0	2	0	0	0	0	5	0	0	0	11	1	0	0		
44,00	0,80	1	27	16	6	3	7	4	3	3	0	2	0	0	0	0	1	0	0	0	5	0	0	0		
82,00	1,49	1	37	22	4	5	8	4	2	8	0	3	0	0	0	0	0	0	0	0	8	0	0	0		
96,00	2,39	4	66	29	6	4	14	5	3	9	6	2	0	0	0	0	0	0	0	0	6	0	0	0		

Сүрөт 11. Кошумча ПЖ чыгарган жыйынтыктардын массиви

Geniş Za	Geniş Za	Gelecek	Gelecek	Gelecek	Gelecek	Şart Kipi	Şart Kipi	Şart Kipi	Emir Kipi	Gerekliiii	Gerekliiii	Gerekliiii	Istek Kip	Istek Kip	Istek Kip	Kullandığı	Toplam	Toplam	Cümle	Sayı
0	0	2	0	0	0	1	2	0	0	0	0	0	0	0	0	103	2			
0	0	5	0	0	0	0	0	0	0	0	0	0	0	0	0	150	17			
0	0	11	0	0	0	0	0	0	0	2	0	0	0	0	0	181	23			
0	0	15	0	0	0	0	0	0	1	1	0	0	0	0	0	179	27			
0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	73	6			
0	0	5	0	0	0	0	0	0	0	0	0	0	0	0	0	106	16			
0	0	5	0	0	0	0	0	0	0	0	0	0	0	0	0	97	13			
0	0	15	0	0	0	1	0	0	2	0	0	0	0	0	0	149	17			
0	0	18	0	0	0	0	0	0	0	0	0	0	0	0	0	93	10			
0	0	10	0	0	0	0	0	0	1	0	0	0	0	0	0	182	19			
0	0	10	0	0	0	0	0	0	0	0	0	0	0	0	0	149	20			
0	0	15	0	0	0	0	0	0	0	0	0	0	0	0	0	110	16			
0	0	10	0	0	0	0	0	0	0	0	0	0	0	0	0	115	12			
0	0	9	0	0	0	0	0	0	0	0	0	0	0	0	0	139	10			
0	0	4	0	0	0	0	0	0	0	1	0	0	0	0	0	124	17			
0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	119	10			
0	0	5	0	0	0	0	0	0	0	4	0	0	0	0	0	181	3			
0	0	14	0	0	0	0	0	0	1	0	0	0	0	0	0	114	16			
0	0	8	0	0	0	0	0	0	1	0	0	0	0	0	0	118	16			
0	0	13	0	0	0	1	0	0	0	0	0	0	0	0	0	259	29			
0	0	5	0	0	0	0	0	0	0	0	0	0	0	0	0	81	11			
0	0	7	0	0	0	0	0	0	0	0	0	0	0	0	0	150	17			
0	0	32	0	0	0	2	0	0	0	0	0	0	0	0	0	243	32			

Сүрөт 12. Кошумча ПЖ чыгарган жыйынтыктардын массивинин уландысы

Машиналык үйрөтүүнү баштоо үчүн алгач дилбаяндар бирдей көптүктөргө бөлүндү. Берилген темага жараша, дилбаянда кездешкен сөздөрдүн санына жана бааларына жараша топторго бөлүп алып, машиналык үйрөтүү жүргүзүлдү. Бардык дил даяндардын жыйынтыктары excel файлында сакталды жана машиналык үйрөтүүнүн жыйынтыктарын сактоо үчүн дагы excel файлы колдонулду. Машиналык үйрөтүүнүн баштапкы этабында, ар түрдүү эксперименттер

аткарылып, алгач берилиш катары кошумча программдык жабдык чыгарып берген жыйынтыктардын бир бөлүгү гана алынып, нейрондук тармак түзүлдү. Алынган жыйынтыктар канааттандыралык болгон жок, башкача айтканда тактык 52% гана түздү. Эксперимент учурунда алынган берилиштер жеткиликтүү эмес экенин байкагандан кийин, кээ бир башка берилиштерди кошуп көрүү менен машиналык үйрөтүү жүргүзүлдү. Изилдөөлөрдүн жыйынтыгында, кээ бир берилиштер баалоо үчүн чоң таасир берип жаткандыгы байкалды. Сүйлөмдүн саны, туура жазылган сөздөрдүн саны, этиштин чактары туура баалоодо эң негизги берилиштердин бир экени такталды. Бул берилиштер менен машиналык үйрөтүү жүргүзгөндө тактык 75% га чейин жогорулады. Ал эми башка берилиштер кошумча, жардамчы берилиш катары колдонулуп, баалоо тактыгын 90% га чейин жакшыртты.

6.2. Баалоо тактыктары

Жадыбал 1 де машиналык үйрөтүү үчүн колдонулган дилбаяндардын көптүгү көрсөтүлгөн. Колдонулган дил баяндардын саны, орточо узундугу, максимум узундугу, минимум баасы, жана максимум бааларына карата ажыратылап, берилиштер топтолду. Топтолгон берилиштер машиналык үйрөтүү үчүн колдонулду.

Жадыбал 1. Машиналык үйрөтүүгө керек болгон дилбаяндардын көптүгү

№	Дилбаян саны	Орточо узундугу	Макс. узундугу	Мин. Баасы	Макс. Баасы
1	1050	350	615	25	92
2	853	350	400	44	89
3	900	150	250	50	95
4	1200	170	345	52	96
5	500	150	423	32	90
6	670	300	476	35	85

Мугалим менен нейрондук тармактын койгон баалары өтө жакын. Демек нейрондук тармак туура үйрөнүп, туура жыйынтык көрсөттү. Нейрондук тармактар 92,5% дык тактыкка чейин туура баалай алды деп айтууга болот. Жадыбал 2 де мугалим аныктаган орточо баа менен нейрондук тармак аныктаган орточо баалардын статистикасы көрсөтүлгөн. Алынган статистиканын жыйынтыгына көз жүгүртсөк, нейрондук тармакты үйрөтүүнүн натыйжасы канааттандырырлык деп эсептесек болот жана үйрөтүлгөн моделди колдонууга болоору далилденди.

Жадыбал 2. Баалоо тактыктары

№	Дилбаян	Мугалим койгон орточо баа	Нейрон койгон орточобаа	Мугалим койгон мин. баа	Нейрон койгон мин. баа	Мугалим койгон макс. баа	Нейрон койгон макс. баа
1	1050	64	54	28	19	96	85
2	853	70	61	32	21	85	79
3	900	55	47	25	15	89	80
4	1200	58	47	24	15	95	88
5	500	65	55	30	20	98	89
6	670	63	54	26	19	95	87

Мугалим жана нейрондук тармактардын баалоолорун салыштыруу аркылуу, машинанын канчалык деңгээлде туура баалоо жүргүзө ала тургандыгын билүүгө болот. Андай моделдердин бири Тау Кендаллдын Каппа Квадраттык Өлчөөсү (QWK- quadratic-weighted kappa) болуп саналат. Бул моделдин жардамы менен машиналык окутуунун сапатын жакшыртууга жана автоматтык баа берүүнү өнүктүрүүгө болот. QWK модели Жадыбал 3тө колдонулуп, жыйынтыктар таблицкага жазылды.

QWK формуласы:

$$W_{ij} = \frac{(i - j)^2}{N - 1}$$

Бул формулада W_{ij} – бул жыйынтык, i – бул мугалим аныктаган баа, j – бул нейрондук тармак аныктаган баа, N – бул колдонулган дилбаяндын саны. Бул модель берген жыйынтыктардан нейрондук тармактын канчалык тактыкта дилбаяндарды баалоого жөндөмдүү экенин билүүгө болот [16].

Жадыбал 3. QWK – моделинен алынган жыйынтыктар

№	Орточо баалардын QWK	Мин. баалардын QWK	Макс. баалардын QWK
1	0,905	0,923	0,885
2	0,905	0,858	0,958
3	0,929	0,889	0,91
4	0,899	0,933	0,959
5	0,800	0,872	0,838
6	0,879	0,927	0,904
AVG	0,886	0,900	0,909

БӨЛҮМ VII.

ТАЛКУЛООЛОР ЖАНА КОРУТУНДУ

7.1 Талкулоолор

Өрнөк баяндамалардын баалоо жыйынтыктарынын негизинде машина үйрөтүү методунун жардамы менен модель курулду. Корутундулоо кадамында, иштелип чыккан модель, програмдык жабдык же плагин түрүнө капсуляцияланат. Табигый тилди изилдөө үчүн алгоритмдерди иштеп чыгуу абдан оор көйгөйлөрдүн бири болуп саналат. Ошондуктан табигый тилде жазылган баяндаманы дагы талдоо жана баалоо иш-чаралары кандайдыр бир тактыкка чейин гана аткарылышы мүмкүн. Албетте, компьютердин жардамы менен баяндама өзгөчөлүктөрүн чыгаруу – бул баяндаманы туура баалоо үчүн жардамчы аспап боло алат, бирок табигый тил татаалдыгы, мугалим сыяктуу толук кандуу баа берүүгө мүмкүнчүлүк бербесин унутпашыбыз керек.

7.2. Корутунду

Бул магистрдикдик диссертация алкагында, окуучулардын баяндама жазуу тапшырмасын баалоо процессин автоматташтыруу, башкача айтканда тапшырма баалоого машина үйрөтүү алгоритми иштелип чыкты. Башкача айтканда, табигый тилди иштетүү куралдарынын жардамы менен окуучулар тарабынан жазылган баяндама же дилбаяны сыяктуу текст тапшырмаларын баалаган системаны иштеп чыгууну жана окуучулар жазган баяндамаларды бир нече параметрлердин негизинде баалаган програмдык жабдык баярдалды. Бул максатка жетүү үчүн, мугалим тарабынан текшерилген баяндамалардын жыйынты топтолуп, баалоо параметрлери иштерип чыкты.

Диссертациялык иштин алкагында иштелип чыккан машина үйрөтүү алгоритминин тактыгы 92,5% түздү. Башкача айтканда, мугалим текшерип

баалаган жана бул иштин алкагында иштелип чыккан нейрондук тармактын койгон баалары өтө жакын болгонун айта алабыз. Демек биз үйрөткөн нейрондук тармак туура түзүлүп, туура жыйынтык көрсөттүп, диссертациялык иштин максатына жетти десек болот. Айткандай эле, нейрондук тармактар 92,5% дык тактыкка чейин туура баалай алды деп айтууга болот. Мисал катары англис тилинде жазылган дилбаяндарды текшерген Educational Testing service I моделин ала турган болсок, бул сервисе колдонулган нейрондук тармактын тактыгы 93-96% түзгөн, башкача айтканда нейрондук тармак дилбаяндардын 93-96% туура баалаган [15].

Жалпысынан бул процессти автоматташтыруу аракеттери көптөн бери изилдөнүүдө. Азыркы учурда актуалдуу болгон массалык ачык онлайн сабактарда колдонуу мүмкүнчүлүгү менен бирге бул түр долбоорлор көптөгөн изилдөөчүлөрдү кызыктырууда. Бирок, табигый тилди изилдөө жана баалоо үчүн алгоритмдерди иштеп чыгуу абдан оор көйгөйлөрдүн бири болуп саналат. Ошондуктан бул диссертациялык иштин алкагында табигый тил иштетүү негиздери даяр болгон түрк тилиндеги баяндамаларды текшерүүнү автоматташтыруу максаты коюлду.

7.3. Алдыдагы изилдөөлөр

Азыркы күндө бул темада көптөгөн илимий изилдөөлөр жүргүзүлүп жатат. Англис тили үчүн жасалган илимий изилдөөлөр жана алгоритмдер көп болгону менен түрк тили жана кыргыз тили үчүн жасалган алгоритмдер дээрлик жокко эсе. Бул изилдөөдө түрк тилинде жазылган дилбаяндарды баалоону автоматташтыруу ишке ашырылды. Бул алгоритмди кыргыз тилинде жазылган дилбаяндар үчүн да колдонууга болот. Ошондуктан кийинки кадам катары кыргыз тили жана башка түрк тилдер үй-бүлөсүндөгү тилдерди үйрөнүү үчүн жазылган дилбаяндарды текшерүү нейрондук тармактар иштелип чыгышы пландалууда.

КОЛДОНУЛГАН БУЛАКТАРДЫН ТИЗМЕСИ:

- [1] Jansen, D., & Schuwer, R. (2015). Institutional MOOC strategies in Europe. Status Report Based on a Mapping Survey Conducted in October-December 2014.
- [2] Barcena, E., Martín-Monje, E., & Read, T. (2015). Potentiating the human dimension in language MOOCs. European MOOCs Stakeholders Summit, 2015 (pp. 6–54).
- [3] Godwin-Jones, R. (2014). Global reach and local practice: The promise of MOOCs. *Language Learning & Technology*, 18(3), 5–15.
- [4] Landauer, T. K. (2003). Automatic essay assessment. *Assessment in education: Principles, policy & practice*, 10(3), 295-308.
- [5] Hirschberg, J., & Manning, C. D. (2015). Advances in natural language processing. *Science*, 349(6245), 261-266.
- [6] Chowdhary, K. (2020). Natural language processing. *Fundamentals of artificial intelligence*, 603-649.
- [7] Gobinda, G. C. (2003). Natural language processing. *Annual Review of Information Science and Technology*, 37, 51-89.
- [8] Mahesh, B. (2020). Machine learning algorithms-a review. *International Journal of Science and Research (IJSR)*. [Internet], 9, 381-386.
- [9] Ayodele, T. O. (2010). Types of machine learning algorithms. *New advances in machine learning*, 3, 19-48.
- [10] Ray, S. (2019, February). A quick review of machine learning algorithms. In *2019 International conference on machine learning, big data, cloud and parallel computing (COMITCon)* (pp. 35-39). IEEE.

- [11] Litjens, G., Kooi, T., Bejnordi, B. E., Setio, A. A. A., Ciompi, F., Ghafoorian, M., ... & Sánchez, C. I. (2017). A survey on deep learning in medical image analysis. *Medical image analysis*, 42, 60-88.
- [12] Hearst, M. A. (2000). The debate on automated essay grading. *IEEE Intelligent Systems and their Applications*, 15(5), 22-37.
- [13] Page, E. B. (1966). The imminence of... grading essays by computer. *The Phi Delta Kappan*, 47(5), 238-243.
- [14] Page, E. B. (2003). Project Essay Grade: PEG. In M. D. Shermis & J. Burstein (Eds.), *Automated essay scoring: A cross-disciplinary perspective* (p. 43–54). Lawrence Erlbaum Associates Publishers.
- [15] Valenti, S., Neri, F., & Cucchiarelli, A. (2003). An overview of current research on automated essay grading. *Journal of Information Technology Education: Research*, 2(1), 319-330.
- [16] Taghipour, K., & Ng, H. T. (2016, November). A neural approach to automated essay scoring. In *Proceedings of the 2016 conference on empirical methods in natural language processing* (pp. 1882-1891).
- [17] Song, S., & Zhao, J. (2013). *Automated essay scoring using machine learning*. Stanford University.
- [18] Lilja, M. (2018). Automatic essay scoring of Swedish essays using neural networks.
- [19] Sharma, S., Sharma, S., & Athaiya, A. (2017). Activation functions in neural networks. *towards data science*, 6(12), 310-316.

[20] Ghosh, S., & Fatima, S. S. (2008, November). Design of an Automated Essay Grading (AEG) system in Indian context. In TENCON 2008-2008 IEEE Region 10 Conference (pp. 1-6). IEEE.



