

әл-Фараби атындағы Қазақ ұлттық университеті

ӘОЖ 004.056

Қолжазба құқығында

**БОЛАТБЕК МИЛАНА АСЛАНБЕКҚЫЗЫ**

**Мәтіндегі экстремистік бағытты анықтау үшін веб-ресурстардағы  
семантикалық талдау модельдерін әзірлеу және зерттеу**

6D100200 – Ақпараттық қауіпсіздік жүйелері

Философия докторы (PhD)  
дәрежесін алу үшін дайындалған диссертация

Отандық ғылыми кеңесші:  
физика-математика ғылымдарының кандидаты,  
қауымдастырылған профессор Мусиралиева Ш.Ж.,  
әл-Фараби атындағы Қазақ ұлттық  
университеті, Алматы, Қазақстан

Шетелдік ғылыми кеңесші:  
профессор Dieter Gollmann  
Гамбург технологиялық университеті, Германия

Қазақстан Республикасы  
Алматы, 2022

## МАЗМҰНЫ

НОРМАТИВТІ СІЛТЕМЕЛЕР .....	4
БЕЛГІЛЕУЛЕР МЕН ҚЫСҚАРТУЛАР .....	5
КІРІСПЕ .....	6
1 ЭКСТРЕМИСТІК МӘЛІМЕТТЕР ТҮСІНІГІ ЖӘНЕ ЕСЕПТІҢ ҚОЙЫЛЫМЫ .....	13
1.1 Экстремизм түсінігі .....	13
1.2 Экстремизмнің жіктелуі .....	14
1.3 Жаһандық экстремистік ресурстар .....	15
1.4 Экстремизмге қарсы күрес .....	19
1.4.1 Экстремизммен күресуге арналған жобалар .....	20
2 ВЕБ-РЕСУРСТАРДАҒЫ ЭКСТРЕМИСТІК МӘЛІМЕТТЕРДІ АНЫҚТАУДЫҢ СЕМАНТИКАЛЫҚ МОДЕЛІН ҚҰРУҒА ҚАЖЕТТІ КОРПУС ҚҰРУ .....	26
2.1 Экстремистік бағытты (ЭБ) анықтау үшін веб-контентті жинауға және талдауға арналған бағдарламалық жабдықтама құрастыру .....	26
2.1.1 Мәліметтерді алуға дайындық кезеңі .....	26
2.1.2 Мәліметтерді жинау .....	28
2.1.3 Мәліметтерді жинау үшін парсер құрастыру .....	30
2.2 Веб-ресурстардағы экстремистік мәліметтерді анықтаудың семантикалық моделін құруға қажетті корпус құру .....	30
2.3 Корпус талдауы .....	32
2.4 Морфологиялық талдау жасауға арналған қосымшаны пайдалану, кілттік сөздерді анықтау .....	33
3 ВЕБ-РЕСУРСТАРДАҒЫ ЭКСТРЕМИСТІК МӘТІНДЕРДІ АНЫҚТАУҒА АРНАЛҒАН СЕМАНТИКАЛЫҚ МОДЕЛЬДІ ЗЕРТТЕУ ЖӘНЕ ҚҰРУ .....	37
3.1 Мәліметтерді алу және препроцессинг модулі .....	38
3.2 Белгілерді шығару және мәліметтерді белгілеу модулі .....	42
3.2.1 TF-IDF .....	43
3.2.2 n-грамдар .....	44
3.2.3 Bag-of-words .....	45
3.2.4 Word2vec .....	46
3.3 Терең оқыту үлгілері арқылы талдау жүргізу модулі .....	47
3.3.1 Терең оқыту үлгілері .....	47
3.3.2 Рекуррентті нейрондық желілер .....	47
3.4 LSTM желілері .....	49
3.5 Ұсынылатын модель .....	50
3.5.1 Word Embedding Layer – Сөз ендіру қабаты .....	50
3.5.2 LSTM-Based Representation – Бейнелеу қабаты .....	51

3.5.3 Үлгінің гиперпараметрлері .....	53
3.5.4 Үлгіні бағалау параметрлері .....	53
3.6 Эксперименталды бөлім .....	55
4 МАШИНАЛЫҚ ОҚЫТУ АЛГОРИТМДЕРІНЕ САЛЫСТЫРМАЛЫ ТАЛДАУ .....	67
4.1 Тірек векторлар машинасы .....	68
4.2 Шешім ағашы .....	70
4.3 Кездейсоқ орман .....	73
4.4 k жақын көрші алгоритмі (k-Nearest Neighbors) .....	75
4.5 Аңқау Байес классификаторы .....	76
4.6 Логистикалық регрессия .....	79
4.7 Градиентті бустинг .....	81
ҚОРЫТЫНДЫ .....	85
ПАЙДАЛАНЫЛҒАН ӘДЕБИЕТТЕР ТІЗІМІ .....	86
ҚОСЫМША А .....	96
ҚОСЫМША Ә .....	98
ҚОСЫМША Б .....	105
ҚОСЫМША В .....	108

## НОРМАТИВТІ СІЛТЕМЕЛЕР

Берілген диссертацияда келесі стандарттарға сәйкес сілтемелер қолданылды:

Қазақстан Республикасы МЖМБС 5.04.034-2011. «Қазақстан Республикасы Мемлекеттік жалпыға білім беру стандарты. Жоғары оқу орнынан кейінгі білім. Докторантура». Негізгі ережелер ҚР БҒМ бекітілген. 17.06.2011 ж. №261. Астана, 2011 ж.

«Диссертацияны безендіру нұсқаулығы», Қазақстан Республикасы БҒМ ЖАК 28 қыркүйек, 2004 жыл. №377-3 ж.

ГОСТ 7.32-2001. Ғылыми зерттеу жұмысының есебі. Безендіру ережелері мен құрылымы.

ГОСТ 7.1-2003. Библиографиялық таспа. Библиографиялық сипаттау.

## БЕЛГІЛЕУЛЕР МЕН ҚЫСҚАРТУЛАР

ШЫҰ	–	Шанхай ынтымақтастық ұйымы
ҰҚК	–	Ұлттық қауіпсіздік комитеті
IBM	–	International Business Machines
LSTM	–	Long short-term memory
TF-IDF	–	Term frequency–inverse document frequency
ЖТН	–	Жеке тіркеу нөмірі
GTD	–	Global Terrorism Database
START	–	Study of Terrorism and Responses to Terrorism
PGIS	–	Pinkerton's Global Intelligence Service
CETIS	–	Center for Terrorism and Intelligence Studies
ISVG	–	Institute for the Study of Violent Groups
RAND	–	Research and Development
RDWTI	–	RAND Database of Worldwide Terrorism Incidents
АҚШ	–	Америка Құрама Штаттары
ASEAN	–	Association of South East Asian Nations
OSINT	–	Open Source INTelligence
SOCMINT	–	Social media intelligence
CBRNE	–	Chemical, Biological, Radiological, Nuclear, and Explosives
DVI	–	Disaster Victim Identification
AI	–	Artificial Intelligence
PHP	–	Hypertext Preprocessor
CGI	–	Common Gateway Interface
ASP	–	Active Server Pages
HTTP	–	HyperText Transfer Protocol
API	–	Application Programming Interface
VK API	–	Vkontakte Application Programming Interface
HTTPS	–	HyperText Transfer Protocol Secure
NLP	–	Natural Language Processing
BOW	–	Bag-of-Words
TF	–	Term Frequency
IDF	–	Inverse Document Frequency
RNN	–	Recurrent neural network
ROC	–	Receiver Operating Characteristic
AUC-ROC	–	Area Under the Curve - Receiver Operating Characteristic
TP	–	True Positive
TN	–	True Negative
FP	–	False Positive
FN	–	False Negative
ML	–	Machine Learning
ЭБ	–	Экстремистік бағыт

## КІРІСПЕ

**Зерттеу тақырыбының өзектілігі.** Қазіргі таңда ақпараттық-коммуникациялық Интернет желісі адамзат өмірінің ажырамас бөлігіне айналды. Адамдар «Твиттер», «ВКонтакте», «Facebook» және т.с.с. әлеуметтік желілерді жаһандық байланыс орнату, пікір алмасу, білім алу және т.б. мақсаттарда пайдалануда. Жеке пайдаланушылардың ғана емес, сонымен қатар ақпараттық ұйымдардың да бүкіл әлемдік кеңістікке белсенді қатысуы ұлттық қауіпсіздікті қамтамасыз ету бойынша ақпараттық-коммуникациялық технологиялар дамуының қазіргі тенденцияларына сәйкес келетін іс-шараларды әзірлеу, атап айтқанда экстремизм мен терроризм идеяларының күшеюіне қарсы тұруға қатысты іс-шараларды ұйымдастыру қажеттілігін анықтайды. Экстремистік ұйымдар жаңа ақпараттық технологияларды қолдана отырып желіде топқа жаңа мүшелер тарту, экстремистік іс-әрекеттерді жоспарлау мен орындау, оқыту жұмыстарын жүргізу, әлеуметке қауіпті әрекеттерді басқару мен координациялауда жасырын ақпарат алмасу, экстремистік іс-әрекеттерді орындау үшін қаржыландыру көздерін іздеу, қолданушыларға, соның ішінде жастарға мақсатты түрде идеологиялық насихат жүргізу мақсатында жабық сайттар құру және т.б. әрекеттерді ұтымды орындауда [1, 2]. Мәселе дүниежүзілік сипатқа ие және жаһандық саяси процестің негізгі қатысушыларының бірі ретінде Қазақстан Республикасы үшін де өте өзекті болып табылады.

Интернет ақпараттық кеңістігі әр түрлі ресурстардан тұрады. Сонымен бірге олардың басым көпшілігі бұқаралық ақпарат құралдары болып табылмайды, нәтижесінде бұқаралық ақпарат құралдары туралы заңнаманың нормаларын қолдану мүмкін болмайды. Экстремизм идеяларының таралуына қарсы тұру үшін қазіргі уақытта құқық қорғау органдары қылмыстық заңнаманың «Экстремизмге қарсы іс-қимыл туралы Қазақстан Республикасының 2005 жылғы 18 ақпандағы №31 Заңы» [3] нормасы, «Қазақстан Республикасында діни экстремизм мен терроризмге қарсы іс-қимыл жөніндегі 2018 – 2022 жылдарға арналған мемлекеттік бағдарламаны бекіту туралы», Қазақстан Республикасы Үкіметінің 2018 жылғы 15 наурыздағы № 124 қаулысы [4] пайдаланылады.

Қазіргі таңда Қазақстан аумағында тыйым салынған террористік құрылымдардың ұлттық тізіміне 22 ұйым енгізілген [5].

2020 жылдың қарашасында ҚР президенті Қ.К.Тоқаев ШЫҰ-ға мүше мемлекеттерді сепаратизм, терроризм мен экстремизммен күреске бағытталған Ақпараттық қауіпсіздік орталығын құруға шақырды [6].

Өкінішке орай, соңғы жылдары қазақстандықтар да экстремистік ұйымдардың қатарына қосылуда. Ұлттық қауіпсіздік комитетінің (ҰҚК) дерегі бойынша, Сирия мен Иракта соғыс басталғалы Таяу Шығысқа 800-ге жуық Қазақстан азаматы кеткен. Олардың көбі – балалар. 2018 жылғы шілдедегі жағдай бойынша, Сирия мен Ирактағы халықаралық террористік ұйымдарда қазақстандық 120 ер адам, 250-ден астам әйел және 500 кәмелетке толмаған бала бар [7].

Елбасы Нұрсұлтан Назарбаевтың тапсырмасымен 2019 жыл қаңтар айынан бастап "Жусан" операциясы аясында Сирия мен Ирак аумағындағы азаматтар

қайтарыла бастады. Комитеттің 2019 жылғы қарашадағы жауабында Сирия мен Ирактан Қазақстанға 277 ересек (57 ер адам, 220 әйел адам) мен 547 кәмелетке жасы толмаған баланың қайтарылғаны, Сирияда әлі 90-нан астам азаматтың бары айтылған [8].

Аталған мәселелерді Қазақстан Республикасының ұлттық қауіпсіздігіне төнетін қатер ретінде қарастыруға болады. Интернеттегі экстремистік іс-әрекеттерге қарсы іс-қимыл саласындағы ахуал күрделі болып қала береді, бұл ғылыми зерттеулер жүргізуді және экстремизмнің кез-келген көріністерін анықтауға, алдын алуға және жолын кесуге бағытталған тиімді және уақытылы шаралар кешенін жүзеге асыруды қажет етеді.

Google, Facebook және Twitter алпауыттары интернеттегі террористік мазмұнды жылдам анықтап, жою үшін жасанды интеллект (AI) технологиясын қолдануға уәде берді [9]. IBM-де жоғарыда аталған әлеуметтік желілердегі барлық деректерді талдай алатын Watson әзірлемесі бар [10]. Ресейде Платонның IT-авторы әлеуметтік желілерді бақылау және қауіп-қатерлерді болжау жүйесін құруда [11]. Германия үкіметі террористік актілерден кейін Интернеттегі террористермен күресу үшін ZITiS атты жаңа киберқауіпсіздік бөлімшесі құрылғанын жариялады [12]. Мұндай жүйелер әзірге Қазақстанда жоқ. Осы себепті Интернеттегі веб-ресурстарға талдау жүргізу, экстремистік мазмұндағы мәтіндерді уақытылы анықтауды автоматтандыру экстремизмге қарсы күрес ұйымдары үшін аса өзекті болып табылады.

Әлеуметтік желілердің дамуы зорлық-зомбылықты жақтаушы, экстремизм мен радикализмді насихаттаушы топтардың жылдам таралуына әсер етуде. Микроблог сайттары, әлеуметтік желі топтарындағы экстремистік мазмұнды анықтауға арналған жұмыстар күрделі және дамып келе жатқан зерттеу саласы болып табылады. Экстремизм мәселесі XX ғасырдан бастап бүгінгі күнге дейін отандық және шетелдік ғылыми әдебиеттерде зерттелуде. Бүгінгі таңда экстремизм әр түрлі ғылыми жұмыстардың зерттеу нысаны болып табылады.

Қазіргі таңда веб-ресурстардағы экстремистік мәтіндерді анықтауға қатысты жұмыстардың басым көпшілігі ағылшын тілі үшін жазылған [13-31]. Ғалымдар веб-ресурстардағы экстремистік бағыттарды анықтау үшін машиналық және терең оқыту әдістерін тиімді пайдалануда [32-46]. Соңғы уақытта веб-ресурстарда неміс [47], орыс [48-52], араб [53-57] тілдерінде жазылған экстремистік мазмұнды қамтитын мәтіндерді анықтауға арналған әр түрлі жүйелер құрылу үстінде. [58-65] жұмыстарда әлеуметтік желі мен микроблогтардан анықталған экстремистік мәтіндердің лингвистикалық ерекшеліктерін анықтауға қатысты зерттеулер жүргізілген, сәйкесінше [66, 67] жұмыстар экстремистік мәтіндердің психологиялық тұстарын зерттеуге арналған. [68-72] жұмыстарда веб-ресурстардағы экстремистік мәтіндерді анықтауды сентимент талдау есебі ретінде қарастыру арқылы шешу ұсынылады.

[1, б.2, 73-77] жұмыстарда әлеуметтік желі топтарындағы экстремистік мәтіндерді анықтауға қатысты зерттеулерге шолу келтіріледі. Веб-ресурстардағы экстремистік мәтіндерді стилметриялық құралдар арқылы анықтауға болатындығы [78] жұмыста келтірілген.

Интернеттегі деректер көлемінің экспоненциалды өсуі деректерді жинаудың басқарылатын әдістерін жасау қажеттілігін тудырады. Экстремистік веб-сайттардан

деректерді жинау үшін автоматты веб-шолғышты құруға қатысты зерттеулер [79-83] жұмыстарда қарастырылады.

Сонымен қатар, экстремизмді зерттеуші сарапшылар жаһандық коронавирустық пандемия салдарынан Интернетте көп уақыт отыруға мәжбүр болған азаматтардың экстремистік ұйымдар үшін осал тұс екенін көрсетеді және аталған факторлардың үйлесуі зорлық-зомбылық экстремизмі мен терроризмнің өсуіне әкеледі деп болжам жасайды [84-86].

Қазақ тілі үшін қылмыстық сипаттағы мәтіндерді анықтауға және талдауға арналған бірқатар ғылыми зерттеулер жүргізілген. Атап айтатын болсақ, Шәріпбай А.Ә. бастауымен Л.Н. Гумилев атындағы Еуразия ұлттық университеті жанындағы «Жасанды интеллект» ғылыми-зерттеу институтының ғалымдарының [87], шаруашылық жүргізу құқығындағы республикалық мемлекеттік кәсіпорын «Ақпараттық және есептеуіш технологиялар институты» Мамырбаев Ө.Ж. еңбектерін келтіруге болады. [88] жұмыста террористік қауіптерге байланысты қазақ тіліндегі мәтіндердің көңіл-күйін талдау үшін сөздікті қолдана отырып, ережеге негізделген әдісті сипаттайды. Онда полярлықты талдау әдістеріне шолу, деректер базасындағы кілт сөздердің мазмұны бойынша беттерді талдайтын парсер, қазақ тіліндегі мәтіндерге морфологиялық, синтаксистік және сентименталды талдау ұсынылған. Ақпараттық қауіпсіздікті қамтамасыз ету, мониторинг және қауіптердің алдын алу кезінде әлеуметтік желілердің пайдаланушысын сәйкестендіруге арналған бағдарламалық кешенді әзірлеу үшін пайдаланылатын онтологиялық білім базасы құрастырылған.

[89, 90] жұмыстарда web-контенттің қылмыстық боялған мәтіндік ақпаратының түрлері (киберқылмыс, террористік акт немесе қаржылық алаяқтық) қарастырылады және мәтіндердегі құқыққа қарсы ақпаратты анықтауға мүмкіндік беретін қолданыстағы лингвистикалық талдау технологиялары талданады. Қазақ тіліндегі мәтіндердегі қылмыстық мәнді атаулы коллакцияларды анықтаудың екі кезеңдік әдісі ұсынылады. Қарастырылатын әдіс атаулы сөзтіркестердің әлсіз құрылымдалған мәтіндердегі автоматты танып ерекшеленудің логикалық-лингвистикалық моделін және коллакцияларды анықтау дәлдігін арттыруға арналған сөз тіркестерінің сөздерінің қисындылығын айырудың ықтималды моделін құрайды.

Жоғарыда келтірілген тұжырымдарты ескере отырып, веб-ресурстардағы қазақ тіліндегі экстремистік бағыттағы мәтіндерді анықтауға арналған модельдерді құру тапсырмасы **аса өзекті** деген қорытындыға келуге болады.

**Диссертациялық жұмыстың мақсаты** – веб-ресурстардағы қазақ тіліндегі экстремистік мәтіндерді анықтау үшін семантикалық талдау моделін кешенді зерттеу және құру.

**Зерттеудің міндеттері.** Қойылған мақсатқа қол жеткізу үшін келесі міндеттерді орындау қарастырылады:

1) Веб-ресурстардағы қазақ тіліндегі экстремистік мәтіндерді анықтау үшін машиналық және терең оқыту әдістерін оқытуға және тестілеуге арналған қазақ тіліндегі экстремистік мәтіндер корпусын құру;

2) Веб-ресурстарда қазақ тіліндегі экстремистік мәтіндерді анықтау үшін семантикалық талдау моделін құру;



3) Веб-ресурстарда қазақ тіліндегі экстремистік мәтіндерді анықтауға арналған әдістерді құру және зерттеу;

4) Қазақ тіліндегі экстремистік түйінді сөздердің тізімін құру;

5) Өзірленген модель мен әдістер негізінде веб-ресурстардағы қазақ тіліндегі экстремистік мәтіндерді анықтайтын бағдарламалық жабдықтаманы құру және тестілеу.

**Зерттеу нысаны.** Веб-ресурстардағы экстремистік бағыттағы мәтіндерді анықтау.

**Зерттеу пәні.** Веб-ресурстардағы экстремистік бағыттағы мәтіндерді анықтауға арналған машиналық және терең оқыту әдістері.

**Зерттеу әдісі.** Зерттеу әдісі ретінде машиналық және терең оқыту әдістері, мәтіндерді жіктеу әдістері, табиғи тілді өңдеу әдістері, нейрондық желілер, элеуметтік желілерді талдау әдістері, статистикалық өңдеу әдістері, жүйелік талдау әдістері қолданыс тапты.

**Алынған нәтижелердің ғылыми жаңалығы:**

1) Алғаш рет қазақ тіліндегі экстремистік мәтіндерді анықтау үшін машиналық оқыту әдістерін оқытуға және тестілеуге арналған қазақ тіліндегі экстремистік мәтіндер корпусы құрылды;

2) Алғаш рет қазақ тілінің ерекшеліктерін ескере отырып, LSTM желісінің сөзді ендіру қабатына алдын ала стемминг алгоритмі орындалған биграммдарға TF-IDF әдісін қолданумен ерекшеленетін және экстремистік мәтіндерді анықтау дәлдігін арттыратын семантикалық талдау моделі құрастырылды;

3) Белгілер жиынтығын қалыптастырудың сөзді ендіру әдістері мен n-граммдарды комбинациялауға негізделетін және экстремистік мәтіндерді жіктеудің сапасын арттыратын әдіс құрастырылды;

4) Алғаш рет қазақ тілінде экстремистік түйінді сөздердің тізімі құрылды.

**Жұмыстың теориялық маңыздылығы.** Диссертациялық жұмыстың теориялық маңыздылығы экстремистік іс-әрекеттер мен ұйымдарды анықтау әдістері мен алгоритмдері саласындағы білім жиынтығына негізделген. Алынған іргелі нәтижелерді әлемдік ғылыми қауымдастық пайдалана алады.

**Жұмыстың практикалық маңызыдылығы.** Әдіс, авторлық куәлік түріндегі қолданбалы нәтижелерді ақпараттық қауіпсіздікті, сыни инфрақұрылымды қамтамасыз ету, интернет-экстремизммен күрес жөніндегі уәкілетті органдар пайдалануы мүмкін.

**Қорғауға шығарылатын негізгі тұжырым.** Машиналық оқыту әдістерін оқытуға және тестілеуге арналған қазақ тіліндегі экстремистік мәтіндер корпусы құрастырылды. Веб-ресурстардағы қазақ тіліндегі экстремистік бағыттағы мәтіндерді анықтауға арналған семантикалық талдау моделі құрастырылды. LSTM моделіне алдын ала стемминг алгоритмі орындалған биграммдарға мәтін ендірудің TF-IDF әдісін қолдану арқылы веб-ресурстардағы қазақ тіліндегі экстремистік мәтіндерді анықтау тапсырмасының дәлдігін арттыруға болатындығы анықталды.

**Сенімділік дәрежесі мен апробациялау нәтижелері.** Зерттеудің сенімділігі мен нәтижелерінің негізділігі міндеттерді қоюдың негізделген жауапкершілігімен, критерийлердің және берілген саладағы зерттеулердің жай-күйінің сарапталуымен,

жүргізілген эксперименттер санының көптігімен, сондай-ақ олардың практикаға табысты енгізілуімен қамтамасыз етіліп дәйектеледі. Диссертация нәтижелері төмендегі ғылыми-әдістемелік конференцияларда баяндалып, талқыланды:

1) Mussiraliyeva Sh., Bolatbek M., Omarov B., Bagitova K. Detection Of Extremist Ideation On Social Media Using Machine Learning Techniques // 12th International Conference on Computational Collective Intelligence. – Vietnam, 2020. – P.743-752, [https://doi.org/10.1007/978-3-030-63007-2\\_58](https://doi.org/10.1007/978-3-030-63007-2_58)

2) Mussiraliyeva Sh., Bolatbek M., Omarov B., Medetbek Zh., Baispay G., Ospanov R. On Detecting Online Radicalization and Extremism Using Natural Language Processing // 21st International Arab Conference on Information Technology (ACIT'2020). – Egypt, 2020, DOI: 10.1109/ACIT50332.2020.9300086

3) Mussiraliyeva Sh., Omarov B., Bolatbek M., Ospanov R., Baispay G., Medetbek Zh., Yeltay Zh. Applying Deep Learning for Extremism Detection // International Conference on Advanced Informatics for Computing Research. – Singapore, 2021. – P.597-605. [https://doi.org/10.1007/978-981-16-3660-8\\_56](https://doi.org/10.1007/978-981-16-3660-8_56)

4) Mussiraliyeva Sh., Bolatbek M., Omarov B., Bagitova K., Alimzhanova Zh. Bigram based Deep Neural Network for Extremism Detection in Online User Generated Contents in the Kazakh Language // International Conference on Computational Collective Intelligence. – Greece, 2021. – P.559-570. [https://doi.org/10.1007/978-3-030-88113-9\\_45](https://doi.org/10.1007/978-3-030-88113-9_45)

5) Болатбек М.А., Создание словаря экстремистских слов для казахского языка, Международная научная конференция студентов и молодых ученых «ФАРАБИ ӘЛЕМІ», Казахстан, Алматы, 2018

6) Мусиралиева Ш.Ж., Болатбек М.А., Әлеуметтік желідегі экстремистік мәтіндерді жіктеу дәлдігін грамматикалық қателерді анықтау және түзету арқылы арттыру, Международная научно-практическая конференция "Актуальные проблемы информационной безопасности в Казахстане, 2020

7) Болатбек М.А., Экстремизм түсінігі. Экстремистік мәтіндерді анықтауға арналған белгілер жинағына шолу, Международная научная конференция студентов и молодых ученых «ФАРАБИ ӘЛЕМІ», Казахстан, Алматы, 2020

8) Болатбек М.А., Экстремистік мәтіндерді сентимент талдау арқылы анықтау, Международная научная конференция студентов и молодых ученых «ФАРАБИ ӘЛЕМІ», Казахстан, Алматы, 2020

9) Байдулла А.М., Мусиралиева Ш.Ж., Болатбек М.А. Экстремистік топтарды анықтау және талдау // Матер. Междунар. научн. конф. студентов и молодых ученых «Фараби әлемі». – Алматы: Қазақ университеті, 2021. – С. 74.

10) Мусиралиева Ш.Ж., Болатбек М.А. Веб-ресурстардағы экстремистік мәтіндерді анықтаудың семантикалық үлгілерін құру және зерттеу // Матер. Междунар. научн. конф. студентов и молодых ученых «Фараби әлемі». – Алматы: Қазақ университеті, 2021. – С. 77.

11) Маден М.Т., Мусиралиева Ш.Ж., Болатбек М.А. Онлайн ортада экстремизмнің лингвистикалық маркерлерін анықтау // Матер. Междунар. научн. конф. студентов и молодых ученых «Фараби әлемі». – Алматы: Қазақ университеті, 2021. – С. 101.

12) Шәріпбекова С.Е., Мусиралиева Ш.Ж., Болатбек М.А. Қазақ тіліндегі он қанатты экстремизмді анықтау үшін веб-контентті жинауға арналған бағдарламалық модуль әзірлеу // Матер. Междунар. научн. конф. студентов и молодых ученых «Фараби әлемі». – Алматы: Қазақ университеті, 2021. – С. 118.

13) Ынтықбай Б.Н., Мусиралиева Ш.Ж., Болатбек М.А. Әлеуметтік желілердегі қауіпсіздік пен конфиденциалдықты машиналық оқыту тәсілдерін қолдану арқылы талдау // Матер. Междунар. научн. конф. студентов и молодых ученых «Фараби әлемі». – Алматы: Қазақ университеті, 2021. – С. 119.

14) Мусиралиева Ш.Ж., Омаров Б.С., Болатбек М.А., Жастай Е. Веб-ресурстардағы қазақ тіліндегі экстремисттік сипаттағы мәтіндерді анықтау // Матер. Междунар. научн. конф. в области информационных технологий, посвященной 75-летию профессора У.А.Тукеева. – Алматы, 2021. – С. 98-104.

**Зерттеушінің жеке үлесі.** Ізденуші диссертациялық жұмыстың қойылған міндеттерін шешті. Веб-ресурстардағы қазақ тіліндегі экстремисттік бағыттағы мәтіндерді анықтаудың семантикалық моделі мен әдісі әзірленді. Машиналық оқыту алгоритмдерін оқыту мен тестілеуге арналған қазақ тіліндегі экстремисттік мәтіндер корпусы құрастырылды. Әзірленген модель мен әдістің тиімділігін анықтау мақсатында эксперименттер жүргізілді. Қазақ тіліндегі экстремисттік кілттік сөздер тізімі құрылды.

**Диссертация тақырыбының ғылыми-зерттеу жұмыстарының жоспарларымен байланысы.** Берілген жұмыс Қазақстан Республикасының Цифрлық даму, инновациялар және аэроғарыш өнеркәсібі министрлігінің тапсырысы бойынша «Мәтіндегі экстремисттік бағытты анықтау үшін веб-ресурстардағы семантикалық талдау модельдерін, алгоритмдерін құрастыру және кибер-криминалистика құрал-жабдықтарын әзірлеу» жобасы аясында жазылды, ЖТН АР06851248.

**Басым бағыт:** Ұлттық қауіпсіздік және қорғаныс

**Мамандандырылған бағыт:** Ақпараттық қауіпсіздікті қамтамасыз ету

**Нәтиженің жарияланымдары.** Диссертация тақырыбы бойынша алынған нәтижелер бес баспалық жұмыста жарияланды. Оның ішінде халықаралық реферативтік мәліметтер қорына енетін басылымдарда жарияланған мақала:

1) Mussiraliyeva Sh., Omarov B., Yoo P., Bolatbek M. Applying Machine Learning Techniques for Religious Extremism Detection on Online User Contents // CMC – Computers, Materials & Continua. – 2021. – Vol. 70. No. 1. P. 915–934. ISSN:1546-2226 (SJR 2020 0.79, процентиль 72) <https://doi:10.32604/cmc.2022.019189>

Қазақстан Республикасы Білім және ғылым министрлігі Білім және ғылым саласындағы бақылау комитеті ұсынатын журналдардағы мақалалар:

1) Bolatbek M.A., Mussiraliyeva Sh.Zh., Tukeyev U.A., Creating the dataset of keywords for detecting an extremist orientation in web-resources in the Kazakh language, Al-Farabi Kazakh National University, Journal of Mathematics, Mechanics and Computer Science, No 1 (97), pp.134-142, 2018

2) М.А. Болатбек, Ш.Ж. Мусиралиева, Экстремисттік мәтіндерді машиналық оқыту әдістері арқылы анықтау, Вестник КазННТУ №6 (130), 300-304 б., 2018

3) К. Шалабаев, К. Әліпбай, М. Болатбек, Ш. Мусиралиева, Вконтакте әлеуметтік желісіндегі экстремистік мәтіндерді анықтау және жіктеу, Вестник КазННТУ №5 (135), 80-86 б., 2019

4) Мусиралиева Ш.Ж., Болатбек М.А., Зият Б.М., Стемминг алгоритмі арқылы экстремистік мәтіндерді жіктеу дәлдігін арттыру, Вестник КазННТУ, №6, 2020 ж.

**Диссертация құрылымы және көлемі.** Диссертациялық жұмыс кіріспеден, 4 бөлімнен, қорытынды, әдебиеттер тізімі және қосымшалардан тұрады. Диссертацияның толық көлемі: 112 бет машиналық мәтіні, оның ішіне 71 сурет, 14 кесте, 137 пайдаланылған әдебиеттер тізімі және 4 қосымша кіреді.

# 1 ЭКСТРЕМИСТІК МӘЛІМЕТТЕР ТҮСІНІГІ ЖӘНЕ ЕСЕПТІҢ ҚОЙЫЛЫМЫ

Бұл бөлімде экстремизм түсінігі, экстремизмнің жіктелуі, жаһандық экстремистік ресурстар, экстремизмге қарсы іс-қимыл, экстремистік іс-әрекеттерге қарсы халықаралық жобалар қарастырылады.

## 1.1 Экстремизм түсінігі

Жаһандық желіде агрессивті ақпараттың таралуына қарсы тұру қоғам мен мемлекеттік органдардың өзекті мәселесі болып табылады, аталған тапсырма интернеттің қажетсіз ресурстарын сүзу арқылы шешіледі.

XXI ғасыр технологиялары интернеттегі ақпаратты пайдалануды кеңейткенімен, мәтіндік деректер интернеттегі мазмұнның ең көп таралған түрі болып табылады. Алайда, террористік және экстремистік топтар ақпарат тарату, насихаттау, қаражат жинау және экстремистік миссияларын қоса алғанда, әртүрлі функцияларды орындау үшін веб-технологияларды ұтымды пайдалануда. Мұндай жағдайда интернет ұлттық қауіпсіздікке қауіп-қатер төндіреді.

Интернет экстремистік материалдарды орналастыру үшін белсенді қолданылып келеді. Проблема жалпы әлемдік сипатқа ие және әлемдік саяси процестің басты қатысушыларының бірі ретінде Қазақстан Республикасы үшін өте өзекті. Интернеттің ғаламдық желісін және компьютерлік байланыс мүмкіндіктерін қолдана отырып, экстремистік қозғалыстар мен топтардың идеологтары азаматтардың, ең алдымен жастардың санасына белсенді әсер етеді.

Экстремистер бүгінде қылмыстық іс-әрекетте Интернет желісінің шын мәнінде шексіз мүмкіндіктерін белсенді пайдаланады, оның ішінде: қылмыстарды дайындау және жасау кезінде Интернет желісінің ресурстарын пайдалану; әлеуметтік қауіпті әрекеттерді басқару және ұйымдастыру мақсатында ақпаратпен жасырын алмасу; экстремистік қызметті қаржыландыру үшін иесіздендірілген қаржылық желілік құралдарды қолдану, мысалы Darkcoin төлем жүйелері; арнайы құрылған сайттарда және басқа интернет-ресурстарда белсенді насихаттау үшін жоспарланған ақпараттық операцияларды жүзеге асыру және т.б. нәтижесінде соңғы жылдары экстремизм проблемасы шиеленісе түсуде, ол қазіргі уақытта жалпы мемлекеттік маңызы бар проблема және Қазақстанның ұлттық қауіпсіздігіне қатер ретінде қаралуда. Осылайша, интернет желісін пайдалана отырып, экстремистік көріністерге қарсы іс-қимыл саласындағы ахуал күрделі болып қалуда, бұл, атап айтқанда, экстремизмнің кез келген көріністерін анықтауға, олардың алдын алуға және жолын кесуге бағытталған ғылыми зерттеулерді жүзеге асыру және тиімді әрі уақтылы шаралар кешенін іске асыру қажеттілігін негіздейді.

Экстремизм құбылысы өте серпінді дамып, күн сайын жаңа белгілер мен сипаттамаларға ие болуда. Қазіргі уақытта үгіт-насихат пен экстремистік идеологияның әсерінен террористердің, сондай-ақ желілік және әлсіз байланысқан құрылымы бар ұйымдасқан террористік қауымдастықтардың, террористік шабуылдарының саны артып келеді. Ақпарат алмасудың және осындай құрылымдарды жылжытудың негізгі құралы – интернет, атап айтқанда веб-ресурстар,

әлеуметтік желілер және электрондық пошта. Осыған байланысты интернет желісінде террористік және экстремистік ақпаратты генерациялайтын және тарататын жекелеген пайдаланушылардан, топтардан және желілік қоғамдастықтардан туындайтын қатерлерді анықтау, қарым-қатынас тақырыптарын, байланыстарды айқындау, сондай-ақ мінез-құлық мониторингі және болжау міндеті туындайды.

Қазіргі уақытта экстремистік әрекетке қарсы күрес Қазақстан Республикасының аумағында да, одан тыс жерлерде де болып жатқан оқиғалармен анықталған құқық қорғау органдарының басым міндеттерінің бірі болып табылады.

Экстремизм – бұл:

– жеке және (немесе) заңды тұлғаның, белгіленген тәртіппен экстремистік деп танылған ұйымдар атынан жеке және (немесе) заңды тұлғалар бірлестігінің әрекеттер жасауы;

– жеке және (немесе) заңды тұлғаның, жеке және (немесе) заңды тұлғалар бірлестігінің мынадай экстремистік мақсаттарды көздейтін әрекеттер жасауы: Қазақстан Республикасының конституциялық құрылысын күштеп өзгерту, егемендігін, оның аумағының тұтастығын, қол сұғылмауын және бөлінбеуін бұзу, мемлекеттің ұлттық қауіпсіздігі мен қорғаныс қабілетіне нұқсан келтіру, билікті күшпен басып алу немесе билікті күшпен ұстап тұру, заңсыз әскерилендірілген құралым құру, оған басшылық ету және қатысу, қарулы бүлік ұйымдастыру және оған қатысу, әлеуметтік, тектік-топтық алауыздықты қоздыру (саяси экстремизм);

– нәсілдік, ұлттық және рулық алауыздықты, оның ішінде зорлық-зомбылықпен немесе зорлық-зомбылыққа шақырумен байланысты алауыздықты қоздыру (ұлттық экстремизм);

– діни өшпенділікті немесе алауыздықты, оның ішінде зорлық-зомбылықпен немесе зорлық-зомбылыққа шақырумен байланысты өшпенділікті немесе алауыздықты қоздыру, сондай-ақ азаматтардың қауіпсіздігіне, өміріне, денсаулығына, имандылығына немесе құқықтары мен бостандықтарына қатер төндіретін кез келген діни практиканы қолдану (діни экстремизм) [3, 91].

## 1.2 Экстремизмнің жіктелуі

Экстремистік көріністердің қолдану саласының сипаты бойынша келесі жіктелуін бөліп көрсетуге болады: саяси сипатта; экономикалық сипатта; діни сипатта және психологиялық сипатта. Кез-келген экстремистік қозғалыс өзара байланысты және әр нақты жағдайда күшті немесе әлсіз көрінетін әртүрлі элементтерді қамтиды. Экстремизмнің негізгі түрлері кесте 1.1-де келтірілген:

Кесте 1.1 – Экстремизмнің негізгі түрлері

Экстремизм түрі	Экстремизмнің берілген түрінің негізгі түсінігі
1	2
Сепаратизм	Мемлекеттің бір бөлігін бөліп алып, оны жаңа тәуелсіз мемлекетке немесе автономды бөлікке айналдыруға ұмтылу
Ксенофобия	Басқа мәдениет, ұлт, мемлекет өкілдеріне деген төзімсіздік

## Кесте 1.1 жалғасы

1	2
Ұлтшылдық	Белгілі бір ұлттың артық екендігін айтуға негізделетін идеология, саяси көзқарастар жүйесі
Шовинизм	Басқа ұлт өкілдерін кемсіту, эксплуатациялау және дискриминациялау мақсатында қандай да бір басқа ұлттық артықшылығын насихаттайтын идеология
Расизм/нәсілшілдік	Әр түрлі нәсілдердің физикалық және психикалық кемелденбегенін насихаттайтын идеология. Халықты «жоғарыдағылар» және «төмендегілер», «кемелденгендер» мен «кемелденбегендер» деген сияқты топтарға жіктеу, нәсілдік тиістілікке байланысты нәсілдік дискриминация, ұлт геноциді үшін пайдаланылады.
Фашизм	Әскери расизм, «басқа» ұлттық және әлеуметтік топтарға ксенофобиядан басталып, геноцидке, мистикалық дұшпандыққа, тоталитарлық мемлекетке табынуға ауысатын әлеуметтік-саяси қозғалыстардың жалпы атауы.
Терроризм	Тек зорлық-зомбылық құралдарын қолдану арқылы орындалатын саясат.

Экстремизмді бағыты бойынша келесідей топтарға жіктеп көрсетуге болады:

1) экономикалық (бір ғана меншік түрін орнату, бәсекелестікті болдырмау және т.б.);

2) рухани (басқа мәдениет өкілдерінің жетістіктерін теріске шығару);

3) экологиялық (табиғатты қорғау саясатына қарсы шығу);

4) діни (басқа конфиссия өкілдеріне деген жеккөрушілік);

5) ұлттық (басқа ұлттардың қызығушылықтары мен құқықтарын теріске шығару);

б) саяси (үкіметтік құрылымдар, мемлекеттік қызметтерге қарсы шығу).

Іс-әрекеттердің масштабына байланысты:

1) мемлекетшілік (өз ұлтына репрессия жасау);

2) мемлекетаралық (өз нормалары мен принциптерін әлемдік масштабға орнатуға тырысу).

Өкілетті құрылымдарға байланысты:

1) мемлекеттік (репрессияның өкілетті құрылымдары арқылы орындалады);

2) мемлекетке оппозициялық (антирежимдік топтар; теракттар).

### 1.3 Жаһандық экстремистік ресурстар

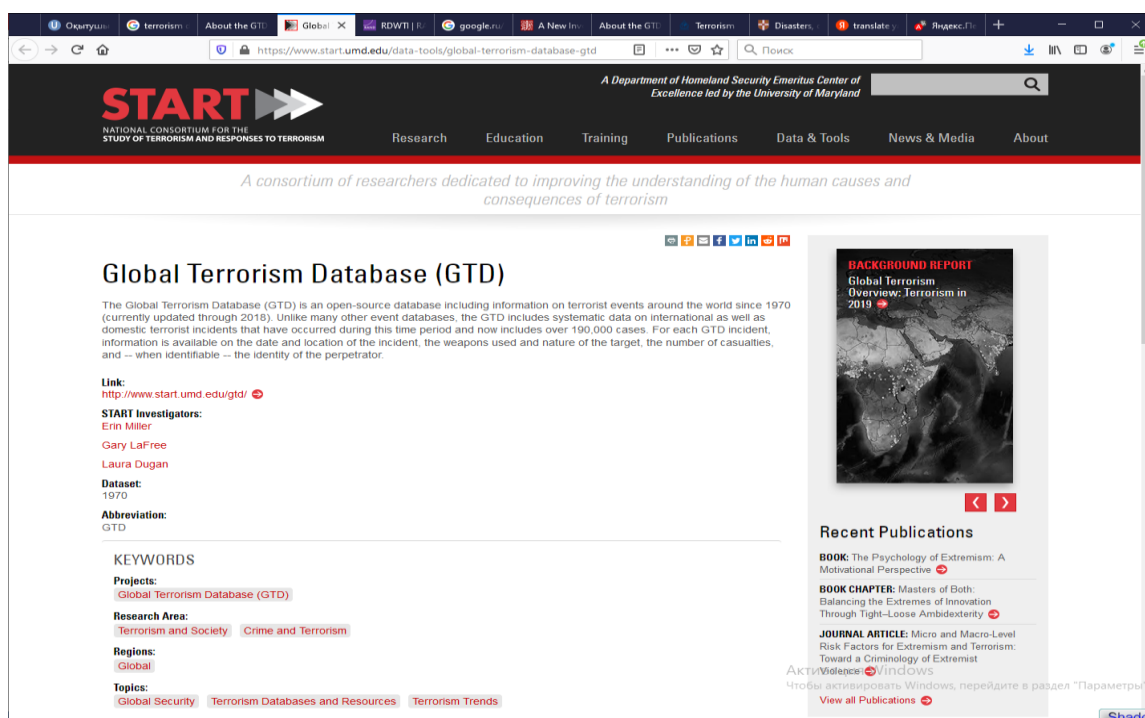
#### *Жаһандық терроризм дерекқоры (GTD)*

Жаһандық терроризм дерекқоры (GTD) – 1970 жылдан 2019 жылға дейін бүкіл әлем бойынша ішкі және халықаралық террористік шабуылдар туралы ақпаратты қамтитын ашық алғашқы дерекқор қоры және қазіргі уақытта 200 000-нан астам жағдайды қамтиды. Әр оқиға үшін оқиғаның күні мен орны, пайдаланылған қару мен мақсаттың сипаты, құрбан болғандар саны және сәйкестендіру кезінде анықталған

топ немесе оған жауапты жеке тұлға туралы ақпаратты қамтиды. Терроризмді зерттеу және терроризмге қарсы іс-қимыл жөніндегі ұлттық консорциум (START) GTD-ді осы онлайн-интерфейс арқылы қол жетімді етеді, террористік зорлық-зомбылық туралы түсінікті кеңейтуге тырысады.

GTD сипаттамалары:

- 200 000-нан астам террористік шабуылдар туралы ақпарат бар;
- Қазіргі уақытта бұл әлемдегі террористік шабуылдар туралы ең толық құпия емес мәліметтер қоры;
- 1970 жылдан бастап 95000-нан астам жарылыс, 20000 кісі өлтіру және 15000 адам ұрлау және барымтаға алу туралы ақпаратты қамтиды;
- Әрбір жағдай үшін кем дегенде 45 айнымалы туралы ақпаратты қамтиды, ал кейінгі оқиғалар 120-дан астам айнымалы туралы ақпаратты қамтиды;
- Тек 1998 жылдан 2019 жылға дейін 4 000 000-нан астам жаңалықтар мақалалары мен 25000 жаңалықтар көздері оқыс оқиғалар туралы мәліметтер жинау үшін талданды (сурет 1.1).



Сурет 1.1 – Жаһандық терроризм дерекқоры (GTD)

Терроризм туралы ғаламдық мәліметтер базасы 2001 жылы доктор Гари Лафри Мэриленд университетіне Пинкертонның Ғаламдық барлау қызметінен (PGIS) қолжазба мұрағатын беруді үйлестірген кезде басталды. 1970 жылдан 1997 жылға дейін PGIS зерттеушілерді негізінен отставкадағы АҚШ әскери күштерін терроризм қаупін бағалау үшін телеграф қызметтерінен, үкіметтік есептерден және ірі халықаралық газеттерден террористік шабуылдар туралы ақпаратты анықтауға және жазуға үйретті. 2005 жылдың желтоқсанына қарай зерттеу тобы түзетулер енгізіп, дерекқорға қосымша ақпарат қосуды аяқтады. 2006 жылдың сәуірінде жаңадан құрылған терроризмді зерттеу және терроризмге қарсы ұлттық консорциум (START),



терроризмді зерттеу және барлау орталығымен (CETIS) бірлесе отырып, 1997 жылдан кейін GTD кеңейту процесін бастады. Бұған жиналған мәліметтер түрлерін кеңейту және pgis-тегі терроризм анықтамасын зерттеу тобы 1970 жылға дейін барлық оқиғаларға ретроактивті түрде қолданған жеке енгізу критерийлерінің жиынтығы ретінде талдау кірді. CETIS командасы 1998 жылдың қаңтарынан 2008 жылдың наурызына дейін болған шабуылдар туралы ақпаратты жазды. Олар сонымен бірге 1993 жылы жетіспейтін деректерді қайта жинауға тырысты. Өкінішке орай, бұл әрекет сәтсіз болды, ішінара сол кезеңнің бастапқы құжаттары жеткіліксіз сақталғанына байланысты. 2008 жылдың сәуірінде Нью-Хейвен университетінің зорлық-зомбылық топтарын зерттеу институтының (ISVG) сарапшылары GTD-ге қосылу үшін мәліметтер жинауды бастады. ISVG-дің 2012 жылдың наурызына дейін жалғасқан мәліметтер жинау әрекеті 2008 жылдың сәуірінен 2011 жылдың қазанына дейін болған террористік шабуылдар туралы мәліметтерді қамтыды. GTD зерттеушілері басында бұл деректерді GTD-ге біріктірді және GTD мүмкіндігінше толық және дәл болуын қамтамасыз ету үшін 1970 жылға дейін бұрын анықталған жағдайлар туралы қосымша істер мен қосымша ақпаратты анықтау үшін көптеген көздерді үнемі қарап отыруды жалғастырды. Аталған процесс GTD құрамына кіретін оқиғаларды анықтау және құжаттау үшін әлемнің әр түрлі бұқаралық ақпарат құралдарынан басталады – күніне екі миллионнан астам мақала қарастырылған. Табиғи тілді өңдеу, атаулы нысандарды алу және машиналық оқыту модельдері террористік шабуылдар туралы ақпаратты қамтитын жаңалықтар мақалаларын анықтауға және ұйымдастыруға көмектеседі. GTD тобы аналитиктерге бірегей шабуылдарды анықтауға, әр оқиғаның егжей-тегжейін жазуға және жаңа ақпарат түскен кезде бұрын тіркелген оқиғалар туралы жазбаларды жаңартуға мүмкіндік беретін жеке деректерді басқару жүйесін жасады. Технологияны жетілдіру және интернетті кеңейту бастапқы материалдардың қол жетімділігін де, жұмыс процестерінің тиімділігін де арттырды [92].

#### *RAND Database of Worldwide Terrorism Incidents*

RAND Database of Worldwide Terrorism Incidents бүкіл әлемдегі терроризмге қатысты оқиғалар туралы RAND дерекқоры (RDWTI) – бұл 1968 жылдан 2009 жылға дейінгі мәліметтер жиынтығы, 40 жыл ішінде RAND корпорациясы терроризм мен терроризмге қарсы зерттеулердің алдыңғы қатарында болды. Осы жұмысты қолдау үшін RAND 1968 жылдан бастап халықаралық және ішкі терроризм туралы жан-жақты ақпарат беретін террористік оқиғалар туралы мәліметтер базасын жасап шығарды. Көптеген жылдар бойы көптеген мемлекеттік және жеке демеушілер RDWTI мен оның алдындағы адамдарға, Rand терроризм хронологиясына және Rand-MIPT терроризм оқиғаларының мәліметтер базасына қолдау көрсетті. Аталған мәліметтер қорында 40 000-нан астам терроризм жағдайлары кодталған және егжей-тегжейлі сипатталған. RAND қызметкерлері аймақтық тәжірибесі, тиісті тілдік дағдылары және елдегі жергілікті жерлерде жұмыс тәжірибесі бар қызметкерлерді тарта отырып, ықтимал террористік шабуылдар туралы кең зерттеу жүргізді. RDWTI – бұл пайдаланушыларға сапалы және жан-жақты мәліметтер беруге арналған толық қол жетімді және интерактивті мәліметтер базасы. Деректер базасы зерттеу және талдау үшін ақысыз және қол жетімді. RAND дерекқоры 1968 жылдан 2009 жылға

дейінгі уақытты қамтиды. Терроризм орындаушылардың жеке басымен немесе себеп сипатымен емес, әрекеттің сипатымен анықталады; негізгі элементтерге мыналар кіреді:

- Зорлық-зомбылық қаупі;
- Қорқыныш пен аландаушылық тудыруға арналған;
- Белгілі бір әрекеттерге мәжбүрлеуге арналған;
- Мотив саяси мақсатты қамтуы керек;
- Әдетте азаматтық мақсаттарға қарсы бағытталған;
- Террористік инциденттер туралы хабарламалардың анықтамалары;
- Инцидент идентификаторы: есеп жасалғаннан кейін әрбір есепке реттік

нөмір беріледі.

Бұл сан RDWTI мәліметтер базасындағы оқиғалардың жалпы санына сәйкес келеді. Оқиға күні: террорлық шабуыл болған күн.

Дереккөз күні: жаңа дереккөздің жарияланған күні.

Ақпарат көзі: ақпарат көзінің атауы. Есептер, әдетте, екі немесе одан да көп дереккөздерге негізделген. Барлық бастапқы құжаттама әр есеп үшін қағаз түрінде сақталады.

Ішкі / халықаралық оқиға: бұл айнымалы үшін әдепкі мән – "ішкі оқиға". "Халықаралық" оқиға деп санау үшін шабуыл элементі шетелдік субъектімен (яғни орындаушы, нысана және т.б.) байланысты болуы керек.

АҚШ-тағы мүлікке/мүлікке шабуыл: әдепкі мән "жоқ"; егер АҚШ азаматы шабуылдың құрбаны болса немесе АҚШ-қа тиесілі мүлік шабуылға ұшыраса немесе бүлінсе, онда "иә" енгізіледі.

Өз-өзіне қол жұмсау миссиясы: әдепкі мән "жоқ"; егер шабуылдаушылар шабуыл аясында "суицид" тактикасын қолданса, онда "иә" енгізіледі.

Шабуыл туралы айтылған: әдепкі мән "жоқ"; егер орындаушылар тобы шабуыл туралы мәлімдесе және бұл мәлімдеме сенімді деп саналса, онда "иә" енгізіледі.

Үйлестірілген шабуыл: әдепкі мән – "жоқ"; егер шабуыл жүйелі түрде жоспарланған жеке шабуылдар сериясының бөлігі болса, онда "иә" енгізіледі.

*Ескертпе:* бір жерде бір мезгілде болған жарылыстар бір инцидент болып саналады.

Үзілген шабуыл: әдепкі бойынша "жоқ" мәні орнатылған; Егер шабуыл шабуылдаушылар жоспарланған шабуылды жүзеге асырмас бұрын тоқтатылса, онда "иә" енгізіледі.

Ел: шабуыл жасалатын ел.

Қала: шабуыл жасалатын қала.

1-Орындаушы: шабуылға жауапты топ. Жауапкершілікті сенімді мәлімдеме арқылы орнатуға болады немесе талдаушы өзінің аймақтық біліміне сүйене отырып, жауапты топқа кіре алады.

2-Орындаушы: егер екінші топ шабуылға жауапты болса, олар осында тізімделеді; әдепкі мән – «жоқ». Жауапкершілікті сенімді мәлімдеме арқылы орнатуға болады немесе талдаушы өзінің аймақтық біліміне сүйене отырып, жауапты топқа кіре алады.

Бірнеше қылмыскер (>2): әдепкі мән «жоқ» [93].

## 1.4 Экстремизмге қарсы күрес

Экстремизмге қарсы іс-қимыл – мемлекеттік органдардың адам мен азаматтың құқықтары мен бостандықтарын, конституциялық құрылыс негіздерін, Қазақстан Республикасының тұтастығы мен ұлттық қауіпсіздігін экстремизмнен қорғауды қамтамасыз етуге, экстремизмнің алдын алуға, оны анықтауға, жолын кесуге және оның салдарларын жоюға, сондай-ақ экстремизмді жүзеге асыруға ықпал ететін себептер мен жағдайларды анықтауға және жоюға бағытталған қызметі.

Экстремизмге қарсы іс-қимыл мынадай негізгі бағыттар бойынша жүзеге асырылады:

- экстремизмнің алдын алуға, оның ішінде оны жүзеге асыруға ықпал ететін себептер мен жағдайларды анықтауға және кейіннен жоюға бағытталған профилактикалық шаралар қабылдау;
- экстремизмді анықтау және оның жолын кесу;
- экстремизмге қарсы іс-қимыл саласындағы халықаралық ынтымақтастық [3].

Мемлекеттік органдар өз құзыреті шегінде экстремизмнің алдын алуға бағытталған мынадай профилактикалық шараларды іске асырады:

1) Діни қызмет саласындағы мемлекеттік реттеуді жүзеге асыратын мемлекеттік орган Қазақстан Республикасының аумағында құрылған діни бірлестіктер мен миссионерлердің қызметіне зерделеу және талдау жүргізеді, өз құзыретіне жататын мәселелер бойынша ақпараттық-насихаттау іс-шараларын жүзеге асырады, Қазақстан Республикасының Діни қызмет және діни бірлестіктер туралы заңнамасын бұзуға қатысты мәселелерді қарайды, діни қызметке тыйым салу туралы ұсыныстар енгізеді;

2) бұқаралық ақпарат құралдары істері жөніндегі уәкілетті орган бұқаралық ақпарат құралдары өнімдерінде экстремизмді насихаттауға және ақтауға жол бермеу, олардың Қазақстан Республикасының заңнамасын сақтауы тұрғысынан мониторинг жүргізеді, мемлекеттік тапсырысты орындайтын бұқаралық ақпарат құралдарында ұлтаралық және конфессияаралық келісімді нығайту мәселелерінің жария етілуін қамтамасыз етеді;

3) білім беру саласындағы орталық атқарушы орган білім беру ұйымдарында білім алушылардың экстремизм идеяларын қабылдамауын, халықаралық құқық пен ізгіліктің жалпыға танылған қағидаттарын құрметтеуін қалыптастыруға бағытталған тәрбиелік бағдарламалардың бекітілуін және іске асырылуын қамтамасыз етеді, оқу орындарының студенттер алмасу мәселелері бойынша халықаралық шарттарының сақталуын бақылауды жүзеге асырады;

4) Қазақстан Республикасының Ұлттық қауіпсіздік органдары жедел-іздігі, қарсы барлау іс-шараларын жүргізеді және Қазақстан Республикасының заңнамасына сәйкес мемлекеттік органдардың дәлелді қорытындылары бойынша өздерінің іс-әрекеттерімен қоғам мен мемлекеттің қауіпсіздігіне қатер төндіретін немесе нұқсан келтіретін шетелдіктер мен азаматтығы жоқ адамдардың Қазақстан Республикасына келуіне жол бермеу жөніндегі шараларды жүзеге асырады;

5) Қазақстан Республикасының Ішкі істер органдары жедел-іздігі қызметін, қоғамдық тәртіпті қорғау және қоғамдық қауіпсіздікті қамтамасыз ету жөніндегі

атқарушылық және өкімдік функцияларды жүзеге асырады, сондай-ақ өз әрекеттерімен қоғам мен мемлекеттің қауіпсіздігіне қатер төндіретін немесе нұқсан келтіретін шетелдіктер мен азаматтығы жоқ адамдарды Қазақстан Республикасының заңнамасына сәйкес Қазақстан Республикасынан шығарып жіберуді жүзеге асырады;

6) облыстардың (республикалық маңызы бар қалалардың, астананың), аудандардың (облыстық маңызы бар қалалардың) жергілікті атқарушы органдары қоғамдық бірлестіктермен өзара іс-қимылды, тиісті аумақтарда құрылған діни бірлестіктер мен миссионерлердің қызметін зерделеуді жүзеге асырады, олар туралы деректер банкін құрады, өздерінің құзыретіне жататын мәселелер бойынша өңірлік деңгейде ақпараттық-насихаттау іс-шараларын жүзеге асырады, Қазақстан Республикасының заңнамасында белгіленген тәртіппен жергілікті атқарушы өңірдегі діни ахуалды зерделеу және талдау тапсырмасын орындайды.

7) Сыртқы барлау субъектілері өз іс-әрекеттерімен қоғам мен мемлекеттің қауіпсіздігіне қатер төндіретін немесе нұқсан келтіретін шет мемлекеттердің ұйымдарына, шетелдіктер мен азаматтығы жоқ адамдарға қатысты Қазақстан Республикасының мемлекеттік органдарын хабардар етуді жүзеге асырады.

Мемлекеттік органдардың экстремизмді анықтау және оның жолын кесу жөніндегі құзыретіне келетін болсақ:

1. Ұлттық қауіпсіздік, ішкі істер органдары және экономикалық тергеу қызметі Қазақстан Республикасының заңнамасымен осы органдардың қарауына жатқызылған қылмыстық құқық бұзушылықтарды анықтайды, жолын кеседі, ашады және тергейді, сондай-ақ Қазақстан Республикасының заңдарында көзделген өзге де өкілеттіктерді жүзеге асырады.

1-1. Экономикалық тергеу қызметі экстремизмді қаржыландыру көздерінің, арналары мен тәсілдерінің алдын алуды, анықтауды, жолын кесуді жүзеге асырады.

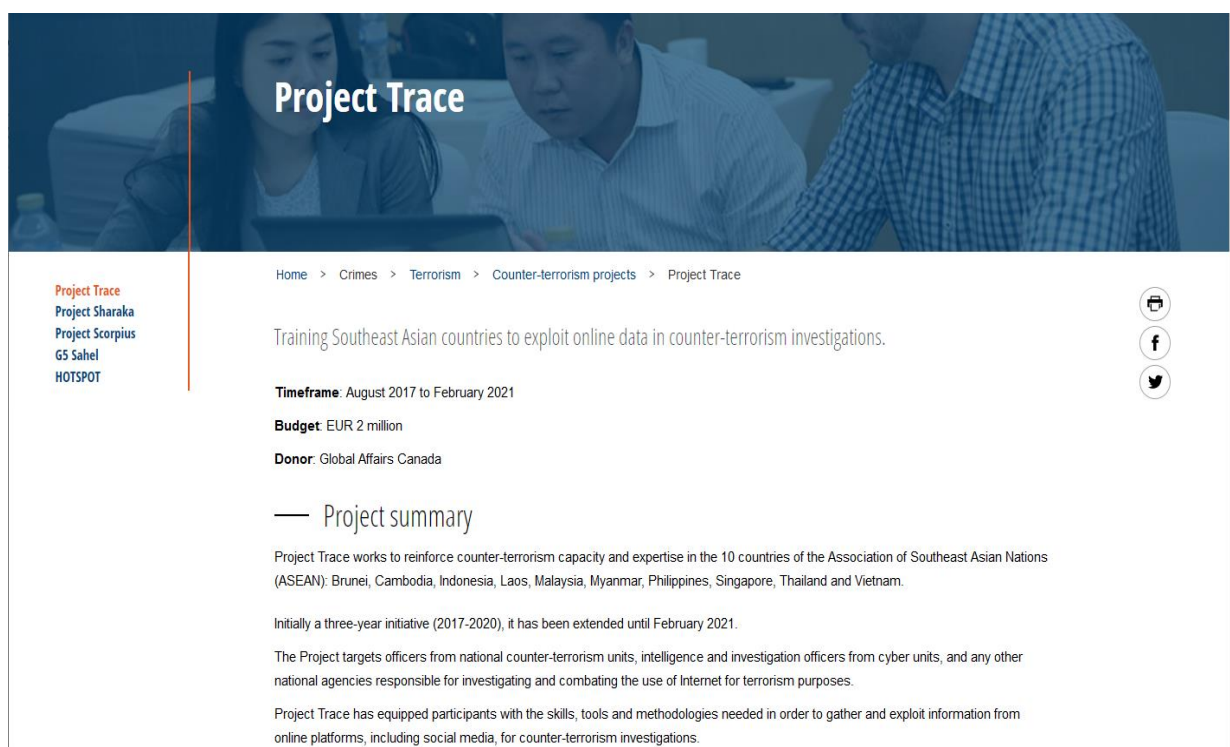
2. Прокурорлар жеке және заңды тұлғалардың, олардың құрылымдық бөлімшелерінің (филиалдары мен өкілдіктерінің) Қазақстан Республикасының Экстремизмге қарсы іс-қимыл саласындағы заңнамасын бұзу фактілері анықталған кезде немесе дайындалып жатқан құқыққа қарсы әрекеттер туралы мәліметтер болған кезде, сондай-ақ бұқаралық ақпарат құралдары арқылы адамның және азаматтың құқықтары мен бостандықтарына, сондай-ақ олардың мүдделеріне зиян келтіруі мүмкін экстремистік материалдар таратылған жағдайда экстремизмнің кез келген көріністерін жою туралы прокурорлық қадағалау актілерін енгізеді, жүзеге асыруға ықпал еткен себептер мен жағдайларды ескере отырып, бұзылған құқықтарды қалпына келтіру туралы сотқа ұйымдар экстремизмді жүзеге асырған жағдайда олардың қызметіне тыйым салу туралы өтініш береді, сондай-ақ Қазақстан Республикасының заңдарында белгіленген тәртіппен және шектерде қылмыстық қудалауды жүзеге асырады.

3. Өзге де мемлекеттік органдар экстремизмді анықтауға және оның жолын кесуге Қазақстан Республикасының заңдарында белгіленген өз құзыреті шегінде қатысады [3, 94].

1.4.1 Экстремизммен күресуге арналған жобалар

*Trace жобасы*

Trace жобасының мақсаты – Оңтүстік-Шығыс Азия елдерін терроризмге қарсы тергеулерде онлайн-деректерді қолдануға үйрету. Мерзімі: 2017 жылдың тамызынан 2021 жылдың ақпанына дейін. Бюджет: 2 миллион евро. Trace жобасы Оңтүстік-Шығыс Азия мемлекеттері қауымдастығының (ASEAN) 10 елінде: Бруней, Индонезия, Камбоджа, Лаос, Малайзия, Мьянма, Филиппин, Сингапур, Таиланд және Вьетнамда терроризмге қарсы күрес саласындағы әлеует пен сараптамалық білімді нығайту бойынша жұмыс істейді. Жоба терроризммен күрес жөніндегі ұлттық бөлімшелердің қызметкерлеріне, кибернетикалық бөлімшелердің барлау және тергеу бөлімшелерінің қызметкерлеріне және интернетті террористік мақсаттарда пайдалануға қарсы тергеу мен күреске жауапты кез келген басқа ұлттық мекемелердің қызметкерлеріне бағытталған. Trace жобасы қатысушыларға терроризмге қарсы тергеу жүргізу үшін онлайн платформалардан, соның ішінде әлеуметтік желілерден ақпарат жинауға және пайдалануға қажетті дағдылар, құралдар мен әдістемелерді ұсынды (сурет 1.2).



Сурет 1.2 – Экстремизмге қарсы Трасе жобасы

Әр цикл келесі әрекеттерді қамтиды:

- Интернетті террористік мақсатта пайдалануға қарсы тұру үшін апталық базалық дайындық;
- Интернетті террористік мақсатта пайдалануға қарсы тұру бойынша бір апталық біліктілікті арттыру курсы;
- Бағдарламалық жасақтаманы сыйға тарту және бағдарламалық жасақтаманы мамандандырылған оқыту;
- Жаттықтырушыға арналған жаттығу;
- Қатысушыларға практикалық тәжірибе беруге арналған үстелдік жаттығу.

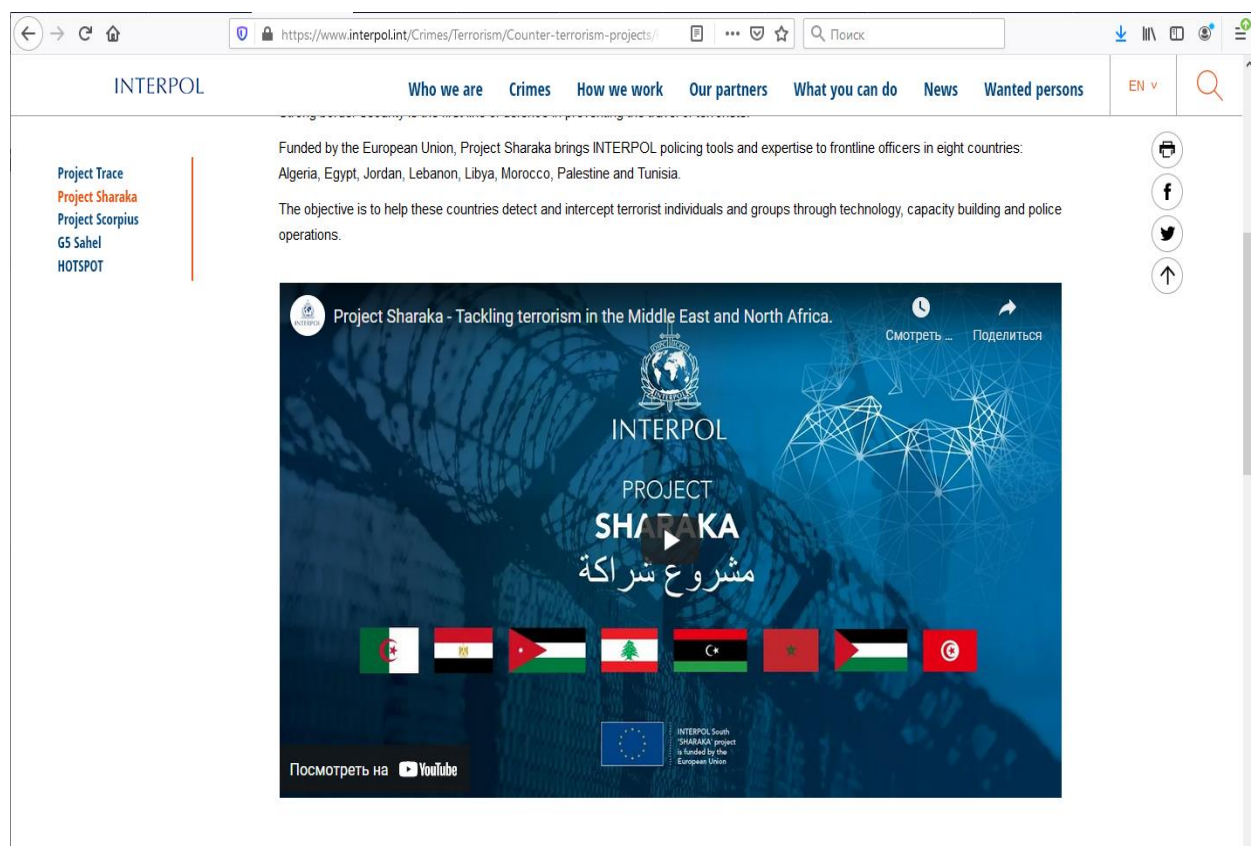
Аталған жобаның мақсаттары келесідей:

- Қылмыстық тергеу аясында ашық бастапқы (OSINT) және әлеуметтік медиа (SOCMINT) деректер ағындарын қалай пайдалануға болатындығын түсіну;
- Интерполдың полиция мүмкіндіктерін, соның ішінде әртүрлі мәліметтер базасы мен интерполдың хабарламаларын қалай пайдалану керектігін түсіну;
- Сандық дәлелдемелерді жинау мен қорғаудың дұрыс әдістерін түсіну;
- Талдауға көмектесетін құрал ретінде графикалық әдістерді қолдану;
- Үшінші тараптардан (Интерпол, басқа да құқық қорғау органдары, интернет-провайдерлер, телекоммуникациялық компаниялар және т.б.) ақпаратты қалай сұрау керектігін түсіну [95].

#### *Sharaka жобасы*

Sharaka жобасы террористердің шекараны кесуінің алдын алуға бағытталған. Аталған жоба озық агенттіктерді I-24/7 (Интерполдың қорғалған жаһандық полиция коммуникациялық желісі), әсіресе әуежайларда, теңіз порттарында және ұлттық шекараларда қосады. Бұл оларға нақты уақыт режимінде барлауды бөлісуге және қылмыс туралы жаһандық мәліметтер базасына қол жеткізуге мүмкіндік береді.

Террористер, әсіресе қақтығыс аймақтарынан оралған шетелдік содыр-террористер пайдаланатын ұрланған жол жүру құжаттарына байланысты шекара қызметі қызметкерлерінің ұрланған және жоғалған жол жүру құжаттары туралы интерполдың деректер базасына тікелей қол жеткізуі өте маңызды (сурет 1.3).



Сурет 1.3 – Экстремизмге қарсы Sharaka жобасы

Бұл жоба мақсатты елдердің терроризмге қарсы күрес саласында қажетті білімге, құрал-жабдықтар мен дағдыларға ие болуын қамтамасыз етеді. Алдыңғы қатардағы қызметкерлер аймақтық тергеулер мен операциялар барысында Интерполдың қылмыстық істер бойынша бірқатар жаһандық дерекқорларын пайдалану бойынша дайындықтан өтеді [96].

#### *Scorpius жобасы*

Scorpius жобасы терроризм мен онымен байланысты трансұлттық қылмыстың алдын алуға және жолын кесуге бағытталған оңтүстік және Оңтүстік-Шығыс Азиядағы құқық қорғау органдарының әлеуетін арттыру жөніндегі екі жылдық (2017-2019 жж.) бастама болды. Оны интерпол мен Канада үкіметі бірлесіп қаржыландырды.

Біртұтас тәсілді қолдана отырып, жоба құқық қорғау қоғамдастығының тиісті қатысушыларын, соның ішінде негізгі шешім қабылдаушыларды, терроризмге және трансұлттық қылмысқа қарсы күрес жөніндегі тергеушілерді, барлау қызметкерлерін, Интерполдың Ұлттық Орталық бюросының қызметкерлерін, қылмыстық сот төрелігі органдарын, прокурорларды және полицияның оқу орындарын біріктірді.

Жаһандық терроризмге қарсы стратегияны қолдай отырып, жоба келесі бенефициар елдерге бағытталған: Бангладеш, Үндістан, Индонезия, Малайзия, Мальдив, Непал, Пәкістан, Филиппин, Шри-Ланка.

Тергеу семинарлары қылмыстық және террористік іс-әрекеттің алдын алу, анықтау және тергеу үшін қажетті практикалық дағдыларды қамтамасыз етеді:

- Күдікті трансұлттық қылмыстық желілер мен олардың филиалдарын анықтау;
- Ұлттық және халықаралық деңгейлерде тиісті құқық қорғау органдары арасында үйлестіру;
- Интерпол және құлақтандыру дерекқорына террористердің бейіндерімен байланысты биометриялық деректерді жүйелі түрде қосу;
- Қару-жарақ пен материалдардың заңсыз айналымын анықтау, қадағалау және жолын кесу;
- CBRNE инциденттеріне және үйдегі жарылғыш құрылғыларға қатысты заттар мен жұмыс әдістері туралы ақпараттармен бөлісу.

Оқу курстары бағалау кезеңінде анықталған қажеттіліктердің әртүрлі салаларына бағытталған, соның ішінде: ашық бастапқы және әлеуметтік медиа тергеулер – бұл ашық бастапқы барлау және әлеуметтік медиа платформаларын қолдана отырып жүргізілген тергеулер арқылы терроризмнің алдын-алудың және онымен күресудің өзекті және тұрақты әдістері. Сондай-ақ, қатысушылар трансұлттық қылмысқа қарсы терроризмге қарсы және онымен байланысты тергеулер барысында жеке өмірге қол сұқпаушылық, адам құқықтары мен сөз бостандығы туралы білді.

CBRNE қару-жарақтары мен материалдары – мүше елдер арасында жаппай қырып-жою қаруы және қолдан жасалған жарылғыш құрылғылар туралы, әсіресе әскери және құқық қорғау органдары арасында барлау алмасуға жәрдемдесу үшін қарастырылады. Сессия сондай-ақ химиялық биологиялық радиологиялық және ядролық және жарылғыш материалдар туындатқан террористік қатерлер мен инциденттерге ден қою жөніндегі елдердің әлеуетін нығайтуға бағытталған.

Қылмыс орнындағы тергеу және апат құрбандарын анықтау (DVI) – жарылыс орнын тексеру, DVI әдістері және терроризмге қарсы күрес контекстіндегі қылмыс орнындағы жалпы тергеу. Қатысушылар мұндай оқиғаларға тиімді жауап беруді және Террористік актілерді жасағандарды қудалауды қолдау үшін өмірлік дәлелдер жинауды үйренді.

Аналитикалық семинарлар офицерлерге қылмыстық және террористік әрекеттерді тиімді талдау үшін қажет практикалық дағдыларды ұсынады:

- Барлау процесінің тұжырымдамалары мен компоненттерін түсіну;
- Талдамалық бағаларды, баяндамалар мен брифингтер дайындау;
- Ашық бастапқы және әлеуметтік медиа зерттеулерінің тұжырымдамаларын, әдістері мен құралдарын түсіну және қолдану;
- Терроризмді қаржыландыру жағдайларын талдау үшін практикалық ақпарат алу [97].

#### *Dark web project AI Lab Dark Web жобасы*

Dark web project AI Lab Dark Web – бұл халықаралық терроризм құбылыстарын (жихадшыларды) зерттеуге және түсінуге бағытталған ұзақ мерзімді ғылыми-зерттеу бағдарламасы. Зерттеушілер халықаралық террористік топтар құрған веб-сайттарды, форумдарды, чаттарды, блогтарды, әлеуметтік желілер сайттарын, бейнелерді, виртуалды әлемді және т. б. қоса алғанда, "бүкіл" веб-контентті жинауға ұмтылады. Олар көп тілді деректерді талдаудың, мәтінді талдаудың және веб-талдаудың әртүрлі әдістерін жасады, сілтемені талдау, мазмұнды талдау, вебкөрсеткіштерді талдау (техникалық күрделілік), көңіл-күйді талдау, авторлықты талдау және зерттеу барысында бейнені талдау есептерін қарастырды. Осы жоба аясында жасалған тәсілдер мен әдістер барлау және қауіпсіздік информатикасы саласын дамытуға ықпал етеді. Мұндай жетістіктер тиісті мүдделі тараптарға терроризм саласында зерттеулер жүргізуге және халықаралық қауіпсіздік пен бейбітшілікті қамтамасыз етуге ықпал етуге көмектеседі. Зерттеушілер сандық кітапханадағы алдыңғы зерттеулер негізінде әртүрлі мамандандырылған өрмекшілер/сканерлер жасады. Құрастырылған өрмекшілер парольмен қорғалған сайттарға кіріп, рандомизацияланған (гуманоидты) үлгіні жасай алады. Олар веб-сайттағы барлық файлдарды, сілтемелерді, PHP, CGI және ASP файлдарын, суреттерді, аудио және бейнелерді шығаруға үйретілген. Жаңалықты қамтамасыз ету үшін әр 2-3 ай сайын таңдалған веб-сайттарды қарап шығады.

Форумдарда өрмекшілерді құруға арналған құрал 15+ форум хостинг бағдарламалық жасақтамасын және олардың форматтарын таниды. Зерттеушілер қатысушылардың өзара әрекеттесуін қайта құруға мүмкіндік беретін авторлар, тақырыптар, жарияланымдар, тақырыптар, уақыт белгілері және т.б. сияқты толық форумды жинайды. Олар форумдардың мерзімді желісін және зерттеу қажеттіліктеріне негізделген біртіндеп жаңартуды жүзеге асырады. Зерттеушілер форумның мазмұнын араб, ағылшын, испан, француз және қытай тілдерінде компьютерлік лингвистиканың таңдаулы әдістерін қолдана отырып жинады және өңдеді.

Көңіл-күйді талдау (полярылық: оң/теріс) және әсер ету (эмоциялар: зорлық-зомбылық, нәсілшілдік, ашу және т.б.) әрі қарай зерттеуді қажет ететін радикалды



және зорлық-зомбылық сайттарын анықтауға мүмкіндік береді. Сондай-ақ, зерттеушілер радикалды идеялардың мазмұнына, жіберушілеріне және олардың өзара әрекеттесуіне негізделе отырып, олардың қаншалықты "жұқпалы" болатындығын зерттейді. Сондай-ақ, олар уақыт пен адамдардың көңіл-күйін/өзгеруін зерттеу үшін жеке визуализация әдістерін жасады. Зерттеуге әсер ететін бірнеше ықтималды көптілді лексика және өлшемді азайту мен болжаудың таңдалған әдістері (мысалы, негізгі компоненттерді талдау) кіреді.

Авторлықты талдау бойынша зерттеулерге негізделе отырып, анонимді жіберушілерді олардың форумдағы хабарламаларына байланысты қолтаңбалар негізінде бірегей сәйкестендіруге мүмкіндік беретін жазу әдістемесін (кибер) әзірледі. Олар дәстүрлі авторлық талдаудың лексикалық және синтаксистік ерекшеліктерін жүйелік (мысалы, қаріп өлшемі, түс, веб-сілтемелер) және семантикалық (мысалы, зорлық-зомбылық) қамтиды [98].

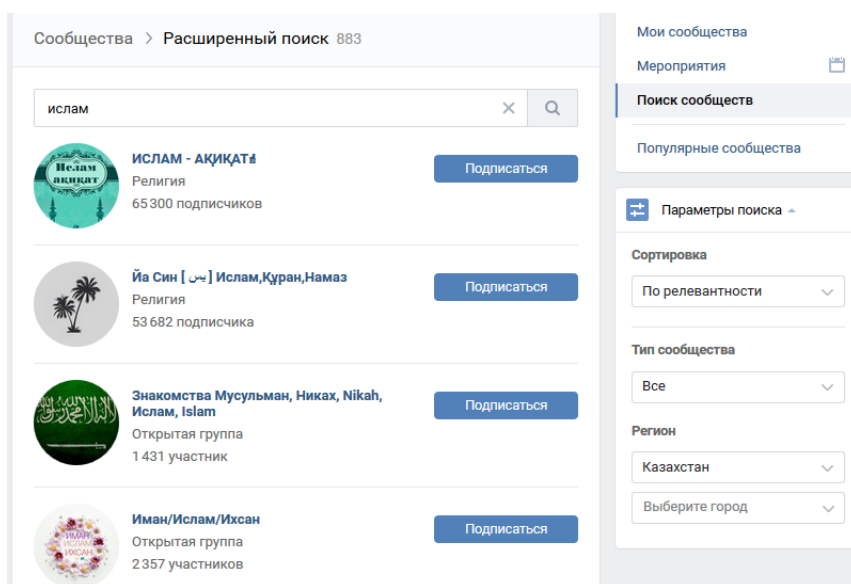
## 2 ВЕБ-РЕСУРСТАРДАҒЫ ЭКСТРЕМИСТІК МӘЛІМЕТТЕРДІ АНЫҚТАУДЫҢ СЕМАНТИКАЛЫҚ МОДЕЛІН ҚҰРУҒА ҚАЖЕТТІ КОРПУС ҚҰРУ

### 2.1 Экстремистік бағытты (ЭБ) анықтау үшін веб-контентті жинауға және талдауға арналған бағдарламалық жабдықтама құрастыру

#### 2.1.1 Мәліметтерді алуға дайындық кезеңі

Веб-ресурстардағы экстремистік мәтіндерді анықтау мақсатында семантикалық талдау модельдерін құру жұмысын жүргізу үшін ең алдымен үлкен көлемдегі корпус қажет болады. Аталған корпус мәтіндері машиналық оқыту алгоритмдерін үйрету және тестілеу кезеңдерінде қолданылады. Бастапқыда ашық қол жетімді ресурстардағы экстремистік мәтіндер іздестірілді. Қазақ тіліндегі экстремистік мәтіндер негізінен жаңалық порталдарындағы, «Youtube», «ВКонтакте» әлеуметтік желілеріндегі комментарийлер арасынан табылды. Алайда мұндай әдіспен табылған экстремистік мәтіндер машиналық оқыту жүйелерін үйрету үшін жеткіліксіз болды. Ашық дереккөздерден шамамен 4400 сөзді қамтитын 400-ге жуық қазақ тіліндегі экстремистік мазмұндағы мәтіндер табылып, корпусқа енгізілді.

Экстремистік мәтіндерді қамтитын корпусты кеңейту үшін қазақ тілді аудитория арасында белсенді қолданылатын әлеуметтік желі түрін таңдау тапсырмасы қойылды. «ВКонтакте», «Facebook», «Twitter» әлеуметтік желілеріндегі мәтіндерге талдау жасалып, нәтижесінде «ВКонтакте» әлеуметтік желісі таңдалды. Келесі кезекте парсер модулінің кірісіне мәтіндері жүктелуі керек болатын топтардың тізімдерін беру қажет. Әлеуметтік желідегі парсинг жасалатын топтар тізімі ашық дереккөздердегі экстремистік мәтіндерге TF-IDF әдісін қолдану арқылы анықталған кілттік сөздер көмегімен құрылды. Атап айтатын болсақ, әлеуметтік желі топтарын анықтау үшін "соғыс", "жиһад", "джихад", "шам", "сирия" және т.б. сөздер қолданылды. "Регион" өрісінде "Қазақстан" таңдалды (сурет 2.1).



Сурет 2.1 – Әлеуметтік желі топтарын іздеу мысалы

Зерттеу жұмысында корпуста экстремистік мәтіндермен қатар бейтарап мазмұндағы мәтіндерді жинау да жоспарланды. Экстремистік мазмұндағы мәтін ретінде экстремистік іс-әрекеттерді орындауға, экстремистік ұйымдарды қаржыландыруға шақыруды, қару-жарақ түрлеріне деген қызығушылықты, қару-жарақ түрлерін қолдан жасауды үйретуді және т.б. қамтитын мәтіндер таңдалды. Бейтарап санаттағы мәтіндерге діни мазмұнды қамтымайтын, қазақстандық аудиторияда кеңінен танымал, жалпы лексикадағы топтар таңдалды.

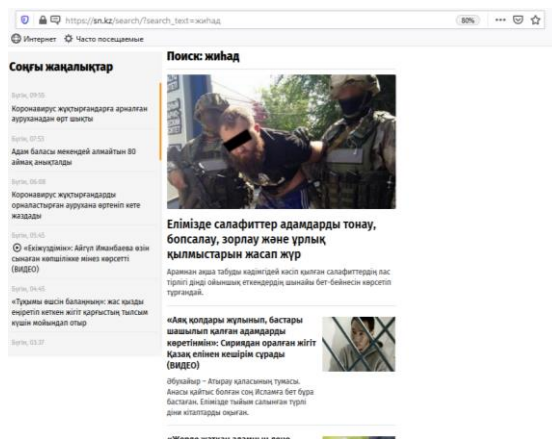
Бірінші кезекте Қазақстан Республикасының территориясында таратуға тыйым салынған топтарды анықтап, оларға парсинг жасау қажет болды. Мұндай топтарды анықтау үшін жоғарыда келтірілген кілттік сөздер қолданылды және табылған нәтижелер ішінен "Берілген материал Қазақстан Республикасының территориясында Қазақстан Республикасының Ақпарат және коммуникация министрлігінің Байланыс, ақпараттандыру және бұқаралық ақпарат құралдары саласындағы мемлекеттік бақылау комитетінің талабы бойынша бұғатталды" деген жазбаны қамтитын топтар таңдалды.

Іздеу нәтижесінде Қазақстан Республикасының территориясында таратуға тыйым салынған 76 топ анықталды. Алайда кей топтарға мүлдем парсинг жасау мүмкін болмады. Ал кей топтардың ішінде экстремистік мазмұндағы мәтіндер табылмады. Қазақстан Республикасының территориясында таратуға тыйым салынған топтар тізімі А қосымшасында келтірілген.

Діни мазмұндағы топтар да жоғарыда сипатталған әдіс бойынша анықталды. Іздеу нәтижесінде 433 топ табылды. Анықталған топтар ішінен мәтіні жоқ, тек суреттерді қамтитын, мәтіні басқа тілдегі, атаулары қайталанған топтар өшірілді. Нәтижесінде корпуста енгізу мақсатында 265 топ таңдалды. Діни мазмұндағы топтар тізімі В қосымшасында келтірілген.

Жалпы лексикадағы топтар да жоғарыдағы әдіс бойынша анықталды. «ВКонтакте» әлеуметтік желісіндегі 100 топ таңдалды. Жалпы лексикадағы топтар тізімі С қосымшасында келтірілген.

Келесі кезекте Қазақстан Республикасында орын алған экстремистік іс-әрекеттер жайлы жаңалық мақалалары жинақталды. Жаңалықтар порталдарға кіріп, "жиһад", "Сирия", "экстремизм" және т.б. сөздерді енгізу арқылы іздестірілді (сурет 2.2).



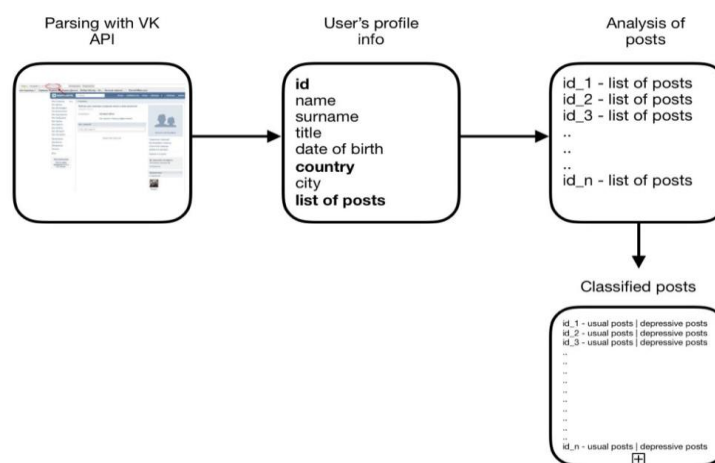
Сурет 2.2 – Экстремистік іс-әрекеттер жайлы жаңалықтарды іздеу процесі

Корпус құрастыру барысында қарастырылған экстремистік іс-әрекеттер жайлы жаңалық мақалаларының тізімі D қосымшасында келтірілген.

### 2.1.2 Мәліметтерді жинау

Ақпаратты экстремистік санатқа жатқызбас бұрын, "қауіп" критерийін анықтау қажет. Шешімдердің бірі – кілт сөздер жиынтығын анықтау. Өзірленген бағдарламалық кешенде ақпарат түрлерін анықтаудың дәл осы әдісі қолданылды.

Анықтау үшін "ВКонтакте" әлеуметтік желісіндегі ақпаратты талдау үшін қолданылатын кілт сөздер жиынтығы жасалды. Бағдарламалық кешен мәтінде көрсетілген кілт сөздердің болуы немесе болмауы негізінде бұл мәтін одан әрі зерттеу үшін жарамды деген қорытынды жасайды [99, 100]. Сурет 2.3-те хабарламалар туралы деректерді жинау, талдау және жіктеудің сұлбасы көрсетілген.



Сурет 2.3 – Деректерді жинау, талдау және жіктеу сұлбасы

Максималды өнімділікке қол жеткізу үшін, ақпарат көздерінен (API) ақпарат алудың кіріктірілген әдістерін қолдану қажет. Егер мұндай әдістер болмаса, HTTP сұраныстары арқылы қажетті ақпаратты алу қажет.

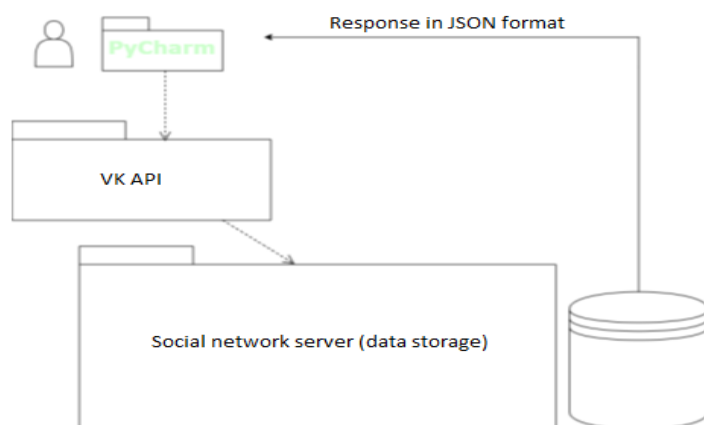
Бағдарламалық кешен үш бөлек модульден тұрады:

1) ақпарат жинау модулі — ашық көздерден ақпаратты қабылдауға және оны әрі қарай өңдеуге беруге жауап береді; «ВКонтакте» әлеуметтік желісінен деректерді талдау үшін Python бағдарламалау тілі қолданылды. Ресми VK API қолдана отырып, қазақстандық профильдер ішінара талданды, сондай-ақ деректер Solr дерекқорында сақталды.

2) кілт сөздерді іздеу модулі – ақпараттың үлкен көлемінің ішінен кілт сөздерді іздеуге жауап береді,

3) құжаттарды саралау модулі – ақпараттың қауіпті екендігін анықтауға жауап береді. Құжаттарды қауіптілік дәрежесі бойынша саралау үшін Long Short Term Memory (LSTM) терең оқыту алгоритмі қолданылды.

Деректерді жинау үшін Тәуелсіз Мемлекеттер Достастығында кеңінен танымал «ВКонтакте» әлеуметтік желісін пайдаланылды. Деректерді жинау үшін Python 3.7 қолданылады [101]. Сурет 2.4-те деректерді жинау процесінің сұлбасы көрсетілген.



Сурет 2.4 – Деректерді жинау сұлбасы

Әлеуметтік желінің API-мен өзара әрекеттесуі сұрау кітапханасы арқылы жүзеге асырылды. Pycharm Community Edition 2018 бағдарламалық жасақтамасы құрастыру ортасы ретінде таңдалды. Деректерді алу үшін «ВКонтакте» API-ді қолданылады – серверге `https` сұрауларын қолдана отырып, «ВКонтакте» әлеуметтік желісінің мәліметтер базасынан қажетті ақпаратты алуға мүмкіндік беретін дайын интерфейс. Кесте 2.1-де пайдаланушының қарапайым сұраныс компоненттері көрсетілген: `'https://api.vk.com/method/users.get?user_id=210700286&v=5.92'`.

Кесте 2.1 – Сұраныс компоненттері

Сұраныс параметрі	Түсіндірмесі
1	2
<code>https://</code>	Қосылу хаттамасы
<code>api.vk.com/method</code>	API қызметінің <code>api.vk.com/method</code> қолданушысы
<code>.get</code>	API Vkontakte әдісінің атауы
<code>?user_id=210700286&amp;v=5.92</code>	Сұраныс компоненті

Әдістер – бұл белгілі бір дерекқор операциясына сәйкес келетін шартты командалар. Мысалы, `пользователи.get` – бұл пайдаланушылар, есептік жазба туралы ақпарат алу әдісі, `.getinfo` ағымдағы пайдаланушы туралы ақпаратты және т. б. қайтарады. Жүйедегі барлық әдістер бөлімдерге бөлінген. Жіберілген сұрауда әдіс атауынан кейін HTTP сұрауында GET параметрлері ретінде кіріс деректерін беру керек. Егер сұраныс сәтті өңделсе, сервер сұралған деректермен JSON нысанын қайтарады. Деректерді талдау үшін `pandas`, `numpy`, `matplotlib`, `plotly`, `bokeh`, `cufflinks`, `sрасу`, `googletrans` пакеттері бар Python 3.7 бағдарламалау тілі есептеу және визуализацияның негізгі кітапханалары ретінде қолданылды [102]. Тиісті мәтіндерді іздеу үшін экстремизммен байланысты кілт сөздер анықталды. Мысалы, "кафир", "өлтіру", "жару" және т.б. бұл кілт сөздер әлеуметтік желілердегі экстремистік посттарды табуға көмектеседі. Экстремистік посттар табылған сайын, кілт сөздер базасы толықтырылып, экстремистік посттардың нақты анықтамасын қамтамасыз етеді. Деректердің аннотациясы үшін «ВКонтакте» әлеуметтік желісінен экстремистік идеялардың мәтіндері жиналды және олардың дұрыс таңбаланғанына көз жеткізу

үшін барлық хабарламалар қолмен тексерілді. Аннотация ережелері мен хабарлама мысалдары кесте 2.2-де келтірілген.

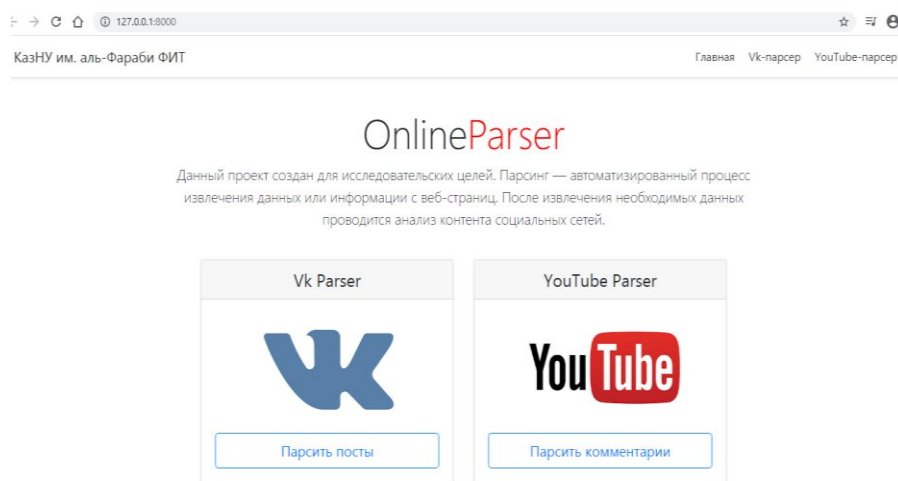
Кесте 2.2 – Аннотация ережесі

Санаттар	Ережелер	Мысалдар
Экстремистік мәтін	(i) Экстремистік ойларды білдіреді.	Мен сыйға ие болу үшін джихад жасау дұрыс деп ойлаймын.
	(ii) Орындалуы мүмкін экстремистік іс-әрекеттерді қамтиды.	Мен Сирия жеріне жиһад жасауға келдім.
Бейтарап мәтін	(iii) Экстремизмге қатысы жоқ.	Маған бұл кітап қатты ұнайды.

### 2.1.3 Мәліметтерді жинау үшін парсер құрастыру

Берілген зерттеу үшін үлкен көлемдегі деректер қажет. Қазақ тілінде дайын экстремистік корпус болмағандықтан, «ВКонтакте», «YouTube», «Твиттер» және елдің жаңалықтар порталдары сияқты ашық ақпарат көздерінен мәлімет жинақтау үшін әр түрлі дереккөздерден ақпарат жинайтын парсер жазылды. Парсердің нәтижесін көрсету үшін шағын веб-қосымша құрылды. Сурет 2.5-те әртүрлі ашық көздерден деректерді жинауға арналған парсер бағдарламасы көрсетілген.

Нәтижесінде, әзірленген парсер көмегімен экстремистік бағыттағы мәліметтер жиналды. Жіктеу модельдерін одан әрі оқыту үшін корпус әзірленді, ол "экстремистік" және "бейтарап" бағыт ретіндегі екі класстан тұрады.



Сурет 2.5 – Зерттеу жұмысы барысында құрастырылған парсер бағдарламалары

## 2.2 Веб-ресурстардағы экстремистік мәліметтерді анықтаудың семантикалық моделін құруға қажетті корпус құру

Корпусқа енгізілген экстремистік сипаттағы хабарламалардың жалпы саны 1200-ге жуық, ал корпустағы жалпы сөздер саны шамамен 140 000 сөзді құрайды.

Құрастырылған корпус .csv форматындағы құжат түрінде сақталған. Экстремистік мәтіндерді қамтитын құжат 4 бағаннан тұрады: хабарламаның реттік нөмірі, хабарлама, хабарламаның препроцессингтен өтіп, тазартылған нұсқасы және қай классқа тиісті жазба екендігін білдіретін "1/0" атрибуты (сурет 2.6).

1.	Лива Туввар Сирия Ал.Лажатта # IslamState-қа қарсы ұрыс бастайды <a href="https://t.co/Fm1iTwk0Z8">https://t.co/Fm1iTwk0Z8</a>	1
2.	террористік иттер бүлікшілерге айналады! #Turkey #PKK #TwitterKurds <a href="https://t.co/...">https:// t...</a>	1
3.	Өл-Хайрдың Газваты 13 арнайы аймақ. Олардан кейін олар Алеппеге қайтып оралады.	1
4.	Сирия-дағы #Ресей   п 'басып кіру' туралы алаңдамаудың кіруі мүмкін емес. Меккедегі бір жұма туралы «Қорғаныс джихад» деп жарияланады ...	1
5.	Ядролық соғыс соғысы келе жатыр !!! <a href="https://t.co/v73GmveabI">https://t.co/v73GmveabI</a>	1
6.	Давлатул ислам! Баакня! Ансарул Халифа - Оңтүстік-Шығыс Азиядағы барлық мухажидтер топтарын жақсы көремін деп сенемін. <a href="https://t.co/KaVLDresnk">https://t.co/KaVLDresnk</a>	1
7.	Орыстар, сириялықтар мен ирандықтар Багдадта әскери үйлестіру камерасын құруда <a href="http://t.co/HYabzUFqZC">http://t.co/HYabzUFqZC</a>	1
8.	@WilayatNinawa сізге хрөсе-ді жауларға тапсырма аласыз!	1
9.	Фаллужаның #IS мандаты Фаллужаның солтүстік-шығысында (12,5) мм қару-жарақпен Сафави әскеріне тиесілі Хамвейде #IRAQ.	1
10.	☐URGENT   #HOMS IS сириялық режимді күшейтеді, 4 сарбазды өлтіруге мүмкіндік береді (Хомс-Пальмира)	1
11.	# Анбар ИСМ басқарылды 2 өлтіру 23 ирак армиясы, Аль-Багдади базасында d әскери жолмен 5 күш жұмсалады	1
12.	Бұл жол # Сириядан емес .. # Фаллужадан #Iraq   мен армияны бомбалау. #ISIS <a href="http://t.co/Qte40bUEI">http://t.co/Qte40bUEI</a>	1
13.	Жана танертен Тел-наам ауданында Нусайристе 10 тонна бомба (халаб) бар.	1
14.	Бұл твит (жүктеме #RISia #ISIS бомбалау емес) төңкерушілерге арналады .. Жана	1

Сурет 2.6 – Экстремистік сипаттағы хабарламалар мысалы

Сурет 2.6-да келтірілген қазақ тіліндегі экстремистік мәтіндерге талдау жасалып, оларға тән кейбір ерекшеліктер анықталды:

1) Қазақ тілінің төл әріптерінің кирилл әріптерімен алмастырылуы (мысалы, «соғысқа» сөзінің орнына «согысқа», «өлтіру» сөзінің орнына «олтиру» және т.б. деп жазу);

2) Биграммалардың жиі кездесуі (мысалы, «жиһад жасау», «Шамға бару», «кәпірлерді қыру» және т.б.);

3) Араб тіліндегі діни терминдердің жиі кездесуі (мысалы, «фард кефайр», «даулатуль ислам» және т.б.);

4) Бір сөздің бірнеше жазылу нұсқасының болуы (мысалы, «джихад», «жиһад», «жихат», «жихад», «джихат»);

5) Корпустағы орфографиялық қателердің көп кездесуі (грамматикалық емес, мәтінді теру барысында пайда болатын типографиялық қателер).

Корпустағы экстремистік мәтіндер келесідей мәліметтерді қамтиды:

– Ауғанстандық "Талибан" қозғалысының ресми баяндамалары ("Жиһад дауысы");

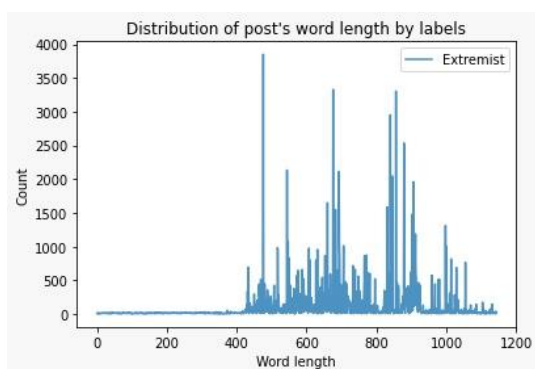
– Аль-Каида ұйымының баяндамалары;

– Экстремистік іс-әрекеттер жайлы қызығушылық танытқан, қазақ тіліне қатысты мәтіндер (жазбалардың қазақша аудармасын сұрау, қазақстандық заңнамаларды көрсету);

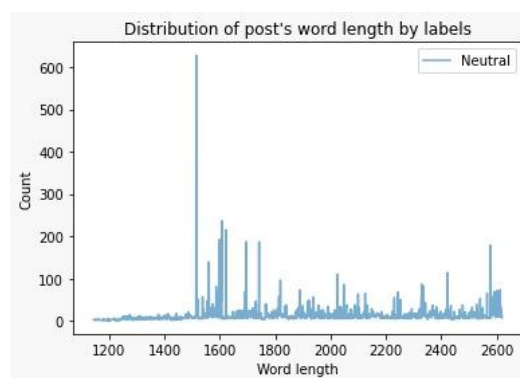
- Құран кітабының жалған аудармасы;
- Қару-жарақ қолдануды үйрету жаттығуларына қатысты жазбалар;
- Есірткіге қатысты жазбалар;
- Жыныстық кемсітушілік, гомосексуализмге қатысты жазбалар;
- Кавказдық жиһадшылардың тізімі көрсетілген жазбалар;
- Сирияға баруға қызығушылық танытқан қолданушылардың жазбалары (Сирияға жету жолдарын сұрау, жиһадтың қалай болып жатқанын сұрау және т.с.с.) [103].

### 2.3 Корпус талдауы

Сурет 2.7-де деректер жиынтығындағы экстремистік және бейтарап мәтіндердің үлестірімі көрсетілген.



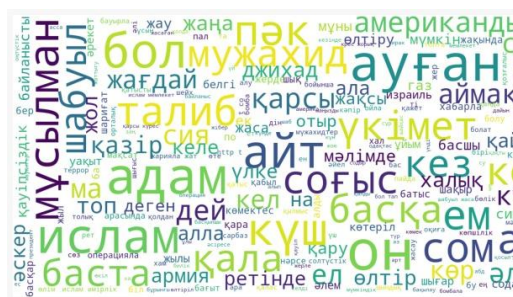
а) Экстремистік хабарламалардың үлестірім графигі



б) Бейтарап мәтіндердің үлестірім графигі

Сурет 2.7 – Корпустағы экстремистік және бейтарап мәтіндердің ұзындықтары бойынша үлестірім кестесі

Сөздер бұлты. Деректерді визуалды түрде көрсету үшін сөз бұлттары қолданылды. Ықтимал экстремистік идеялары бар пайдаланушылардың хабарламалары сурет 2.8-де бөлек көрсетілген.



а) Экстремизмге байланысты хабарламалардың сөздер бұлты



б) Экстремизммен байланысыз хабарламалардың сөздер бұлты

Сурет 2.8 – Экстремистік және бейтарап мәтіндер корпусын визуализациялау



Экстремистік посттар пайдаланушылардың экстремистік ойларына тікелей нұсқау бере отырып, "қарсы" (қарсы), "әскери" (әскери) сияқты сөздерді жиі қолданады [104].

## 2.4 Морфологиялық талдау жасауға арналған қосымшаны пайдалану, кілттік сөздерді анықтау

Табиғи тілді өңдеу саласындағы маңызды мәселелердің бірі – стемминг, яғни кіріс мәтіндегі сөздердің негіздері мен қосымшаларын автоматты түрде анықтау және ажырату. Стемминг кезеңін іске асыруға арналған бірнеше әдістер бар және олардың көбісі бастапқы тілге тәуелді болып келеді.

Кіріс мәтінге талдау жасау барысында стемминг тапсырмасының орындалуы жіктеу дәлдігін айтарлықтай жоғарылатуға септігін тигізуі мүмкін, мысалы стемминг тапсырмасы орындалмаған жағдайда классификациялау жүйесі мысал ретінде мәтіндегі «сабаққа», «сабақтың», «сабақта», «сабақ» сөздерін әр түрлі сөз ретінде таниды, ал аталған сөздердің қосымшаларын қарастырмай, стемминг қадамын орындап, сөз негіздерін алатын болсақ, онда қосымшалардың барлығы алынып тасталатындықтан, негізі бір сөздердің барлығы бір сөз ретінде анықталатын болады. Сондай-ақ, стемминг тапсырмасын орындамаған жағдайда мәліметтер қорына сөздіктегі сөздердің барлық морфологиялық нұсқасын енгізу қажет болады, ал бұл өте үлкен жадыны қажет етеді және сәйкесінше жүйенің жұмысын айтарлықтай баяулатады.

Морфологиялық анализатордың кірісіне қазақ тіліндегі қосымшасы бар сөздер беріледі. Кіріс мәтіндегі сөз негіздерін анықтау үшін қазақ тіліндегі 24 мың және 76 мыңға жуық сөз негіздерін қамтитын екі дерекқор құрастырылды (сурет 2.9). Кестеге экстремистік мәтіндерге тән кілт сөздер де енгізілді.

id	kaz	type
Click here to define a filter		
9928	норма	noun
9929	вахабис	noun
9930	жихад	noun
9931	өлтір	verb
9932	жихад	noun
9933	соғыс	noun
9934	қыр	verb
9935	жихат	noun
9936	джихад	noun
9937	джихат	noun
9938	соғыс	noun
9939	шайқас	noun
9940	шайқас	noun
9941	кафир	noun
9942	кафир	noun

Сурет 2.9 –76000 сөз негізінен тұратын деректер қоры

Стемминг модельді оқыту үшін өзгертілген нысандарды енгізу және үлгі ережелерінің ішкі жиынтығына сәйкес түбірлік нысанды генерациялау арқылы орындалады, ең тиісті ережелерді немесе ережелердің бірізділігін қолданумен, сондай-ақ сөз негіздерін таңдаумен байланысты шешімдер нәтижелік дұрыс сөздің ең жоғары ықтималдығы болуы негізінде қолданылады.

Жалғау — сөз бен сөзді байланыстыратын, сөз аралығындағы қатынастардың көрсеткіші болып табылатын, сөзге грамматикалық мағына үстейтін қосымшалар.

Жалғаудың төрт түрі бар:

- Көптік жалғау;
- Тәуелдік жалғау;
- Септік жалғау;
- Жіктік жалғау.

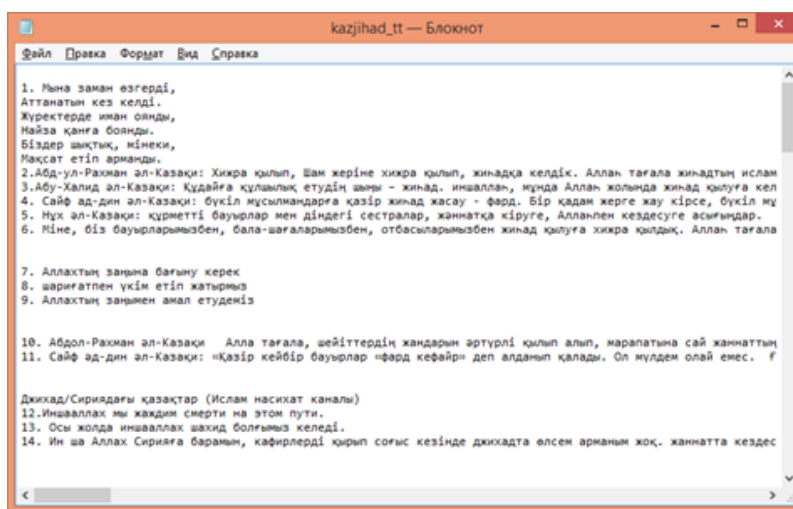
Жалғаулар бірінен соң бірі жалғана береді. Мұндай жағдайда көбінесе алдымен көптік, онан соң тәуелдік, сөз соңында септік жалғаулары жалғанады, жіктік жалғауы да сөз соңында жалғанады: оқушы-лар-ымыз-ға, бала-мыз, келе-сіз [105].

Жұрнақ – жалғанған сөзінен жаңа сөз тудыратын немесе сөзді түрлендіретін қосымша. Қазақ тіліндегі жұрнақ мағынасы мен қызметіне қарай екіге бөлінеді:

1) сөз тудыратын жұрнақтар өзі жалғанған сөзінен жаңа сөз тудырады. Мысалы, “жылқы-шы”, “біл-ім”, “жасы-қ”, “таға-ла”;

2) сөз түрлендіретін жұрнақтар өзі жалғанған сөзіне үстеме мағына қосып, сөздің тұлғасын өзгертеді. Мысалы, “көк-шіл”, “көк(г)-ірек”, “сары-лау”, “сары-рақ”, “жаз-ып”, “жаз-ғалы”. Жұрнақтар сөзге белгілі бір жүйеде рет-ретімен жалғанады [106].

Берілген жұмыста қазақ тіліндегі экстремистік мәтіндерді анықтау дәлдігін арттыру мақсатында қазақ тілі үшін [107] жұмыста ұсынылған стемминг алгоритмі пайдаланылды. Мысал ретінде 4400 сөзден тұратын мәтін қарастырылды (сурет 2.10). Кіріс мәтіні input.txt файлында сақталады. Арнайы жазылған бағдарлама кіріс мәтіндегі сөздердің қосымшаларын бөлу алгоритмін қолданады.



Сурет 2.10 – Кіріс мәтін

Мәтіндегі әр сөздің жұрнағы жоғарыда сипатталған қосымшаларды бөлу алгоритміне сәйкес алынып тасталады, тек негіздер қалады. Кіріс мәтіндегі сөздердің негіздерін анықтау үшін қазақ тіліндегі 24 мың және қазақ тіліндегі 76 мыңға жуық сөздік негіздерден тұратын екі дерекқор пайдаланылды (сурет 2.11).



Сурет 2.11– Бағдарлама нәтижесі

Сурет 2.12-де жоғарыда сипатталған қосымшаларды бөлу алгоритмімен анықталған кіріс мәтініндегі сөздердің негіздері көрсетілген. Бағдарлама жоғарыда ұсынылған алгоритм бойынша мәтіннен тек сөздерді іріктеп алып, олардың түбірін таба тауып, output.csv файлына жазып отырады.

3	Кіріс сөз	Сөз негізі	28	хижа	хижа
4	манна	манна	29	қыпып	қып
5	заман	заман	30	ш ам	ш ам
6	өзгерді	өзге	31	жеріне	жер
7	атпанын	атпан	32	хижа	хижа
8	кез	кез	33	қыпып	қып
9	келді	кел	34	жиһадқа	жиһад
10	жүректерде	жүрек	35	келдік	кел
11	иман	иман	36	аллаһ	алла
12	оянды	оян	37	тағала	таға
13	найза	найза	38	жиһадтың	жиһад
14	қанға	қан	39	исламның	ислам
15	баянды	бая	40	ең	ең
16	біздер	біз	41	биігі	биік
17	шықтық	шық	42	екенін	екенін
18	мінеки	міне	43	көрсеті	көрсет

Сурет 2.12 – Бағдарламаның нәтижелері жазылған .csv форматындағы файл

Сөздің дұрыс негізін табу дәлдігі мәліметтер базасындағы сөздердің санына байланысты. Деректер базасындағы сөздердің аз саны дәлдікті төмендетеді. Осы себепті кілт сөздер базасын экстремистік сипаттағы сөздермен толықтыру қажет. Сондай-ақ, кейбір жұрнақтар сөздің негізін өзгерте алатындығын ескерген жөн. Мысалы, «халық» сөзіне «ның» қосымшасын қосқан кезде «халқының» сөзін аламыз.

Стемминг алгоритмін қолдана отырып, машиналық оқыту әдістері көмегімен қазақ мәтіндерін экстремистік және бейтарап топтарға жіктеу міндеті орындалды [108, 109].

TF-IDF (ағылшынша TF — term frequency, IDF— inverse document frequency) — құжаттар жиынтығының немесе корпустың бөлігі болып табылатын құжат контекстіндегі сөздің маңыздылығын бағалау үшін қолданылатын статистикалық шама. Кейбір сөздердің салмағы құжаттағы осы сөзді қолдану жиілігіне пропорционал және коллекцияның барлық құжаттарындағы сөзді қолдану жиілігіне кері пропорционал. TF-IDF өлшемі көбінесе мәтінді талдау және ақпаратты іздеу тапсырмаларында, мысалы, кластерлеу кезінде, құжаттардың жақындық өлшемін есептеу кезінде, құжаттың іздеу сұранысына сәйкестігі критерийлерінің бірі ретінде қолданылады. TF-IDF-те белгілі бір құжаттың ішінде жиілігі жоғары және басқа құжаттарда қолдану жиілігі төмен сөздер жоғары мәнге ие болады. TF-IDF әдісі қазақ тілінде құрастырылған корпустық мәтіндердегі кілт сөздерді анықтау үшін қолданылды. Салынған корпуста барлық сөздердің TF-IDF мәндері анықталды. Сурет 2.13-те діни сипаттағы сөздер және олардың TF-IDF мәндері көрсетілген [110, 111].

1	words, TF-IDF
2	алла, 1064.94225986155
3	намаз, 374.3576272057989
4	пайғамбар, 351.0525812347899
5	раббы, 261.4808053326351
6	мұсылман, 229.10748607400885
7	құран, 220.06105622628598
8	аят, 218.08389059261924
9	иман, 208.80295300866857
10	мұхаммед, 199.28104933233564
11	дұға, 181.22857579929922
12	тағала, 178.858674885299
13	ораза, 173.9963381928556
14	дін, 173.26824843258433
15	елші, 172.20982415978003
16	хадис, 171.07026346804878
17	ислам, 160.7247364573076
18	күнә, 157.92112765006047
19	ибн, 145.30556994837147
20	қан, 122.86302521944057
21	тозақ, 116.32048363840886
22	қиямет, 114.50660510201877
23	имам, 112.11613282128268
24	құлшылық, 109.9706927331165
25	сауап, 106.84141213233139
26	періште, 105.84399921405313
27	азап, 102.65485548483262
28	күш, 102.444540432683

Сурет 2.13 –TF-IDF мәні бойынша сұрыпталған сөздер

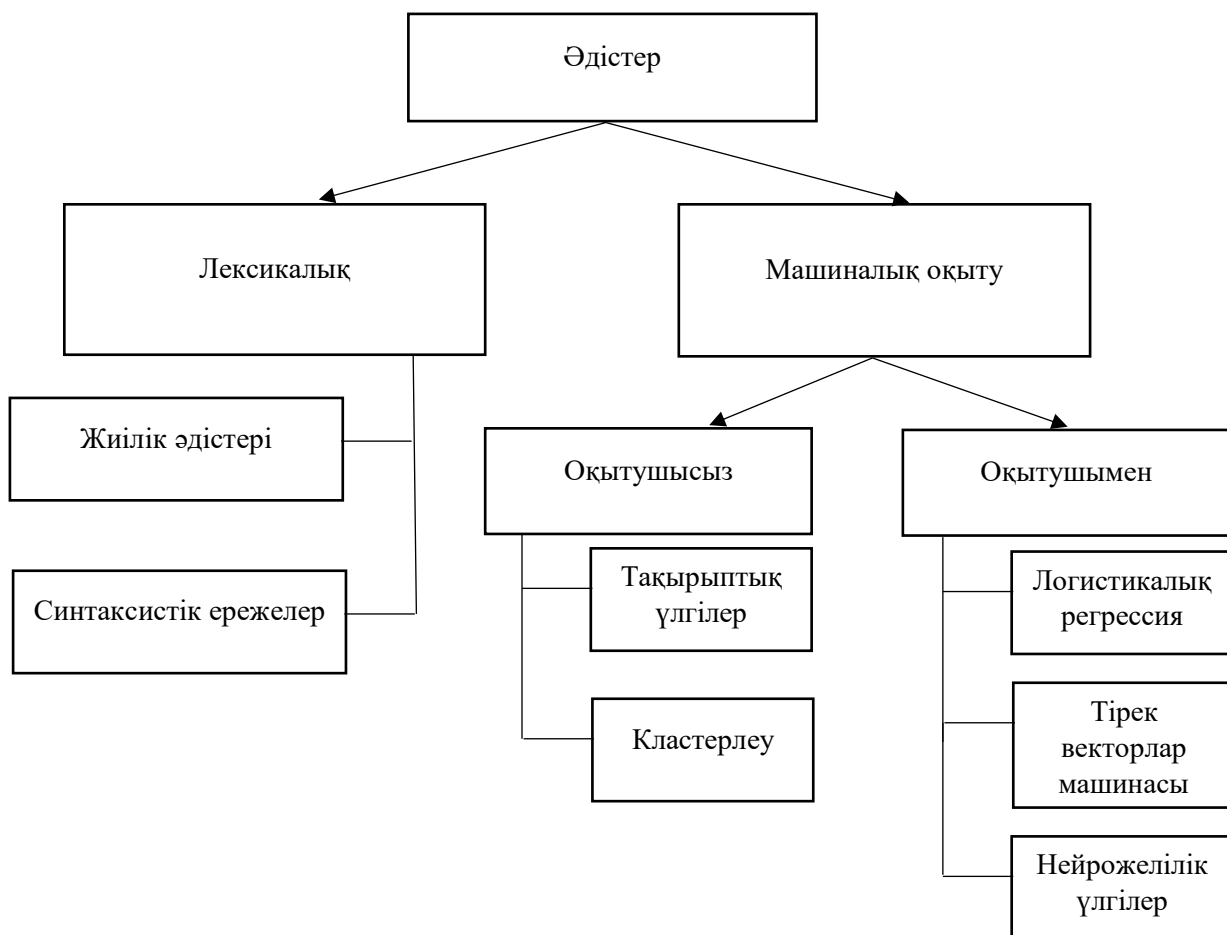
Аталған діни сипаттағы сөздер таңдалды және олар жаңа құжатқа кіріс мәтінді машиналық оқыту әдістерімен жіктеуде маркер ретінде пайдалану үшін сақталды.

### 3 ВЕБ-РЕСУРСТАРДАҒЫ ЭКСТРЕМИСТІК МӘТІНДЕРДІ АНЫҚТАУҒА АРНАЛҒАН СЕМАНТИКАЛЫҚ МОДЕЛЬДІ ЗЕРТТЕУ ЖӘНЕ ҚҰРУ

Мәтінді жіктеу – бұл әрбір  $(d_j, c_i) \in D \times C$  жұбына логикалық мәнді тағайындау есебі, мұндағы  $D$  – құжаттар аймағы, ал  $C = \{c_1, \dots, c_{|C|}\}$  – алдын ала анықталған санаттар жиынтығы.  $(d_j, c_i)$ -ге тағайындалатын  $T$  мәні  $d_j$  құжатын  $c_i$ -ге жатқызу туралы шешімді, ал  $F$  мәні белгілі бір жағдайларға байланысты  $d_j$  құжатын  $c_i$ -ге жатқызбау туралы шешімді білдіреді. Есептің ресми түрдегі қойылымы белгісіз мақсатты  $\dot{\Phi}: D \times C \rightarrow \{T, F\}$  функциясын (құжаттарды қалай сипаттау керектігін сипаттайтын) классификатор деп аталатын  $\Phi: D \times C \rightarrow \{T, F\}$  функциясы арқылы (ереже, гипотеза немесе модель деп аталуы мүмкін)  $\dot{\Phi}$  пен  $\Phi$  «максималды түрде сәйкес келетіндей етіп» жуықтау болып табылады.

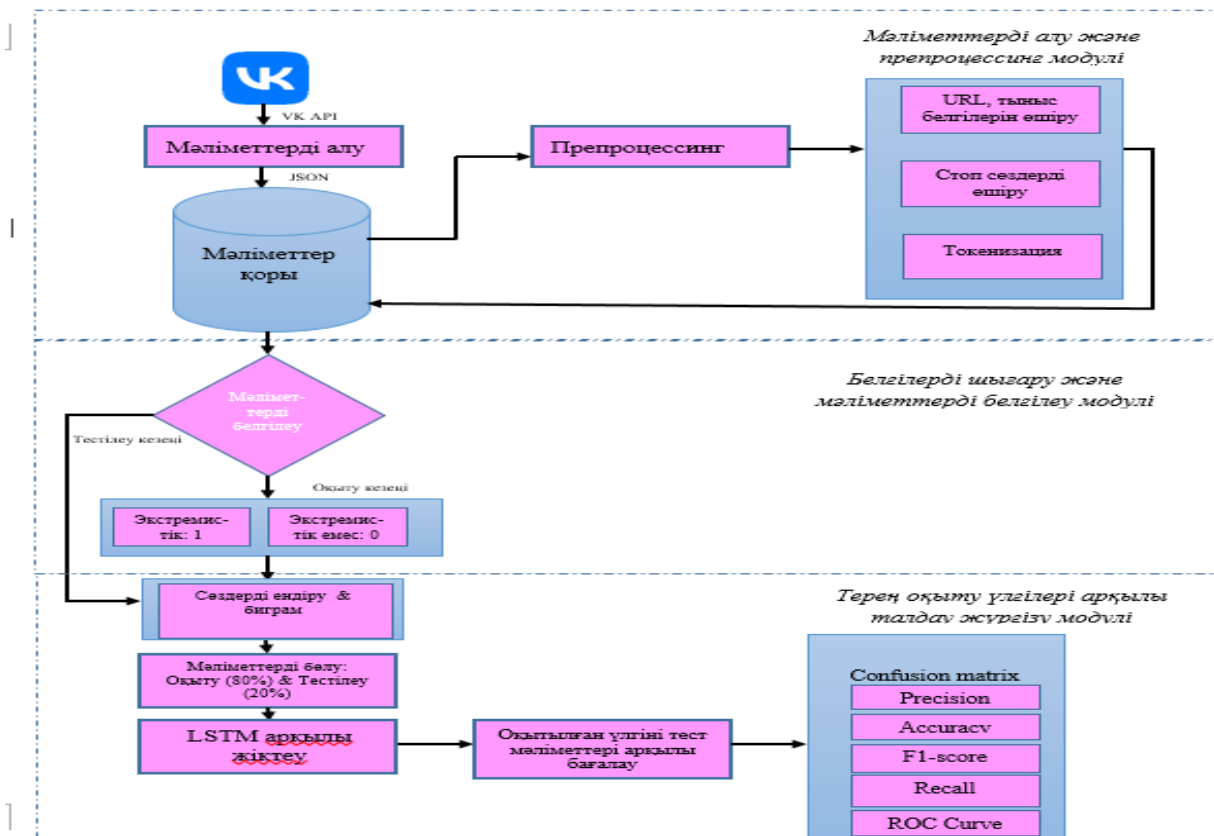
$C = \{c_1, \dots, c_{|C|}\}$  бойынша жіктеу  $D$ -дағы құжаттарды берілген  $c_i$ , мұндағы  $i = 1, \dots, |C|$  санатына жіктеудің  $|C|$  тәуелсіз есебі ретінде қарастырылады.  $C_i$  үшін классификатор деп белгісіз  $\Phi_i: D \rightarrow \{T, F\}$  мақсатты функциясын жуықтайтын  $\dot{\Phi}_i: D \rightarrow \{T, F\}$  функциясын атаймыз [112].

Мәтінді жіктеу есебін шешуге арналған әдістерді келесідей топтарға бөлуге болады (сурет 3.1):



Сурет 3.1 – мәтінді жіктеу әдістері

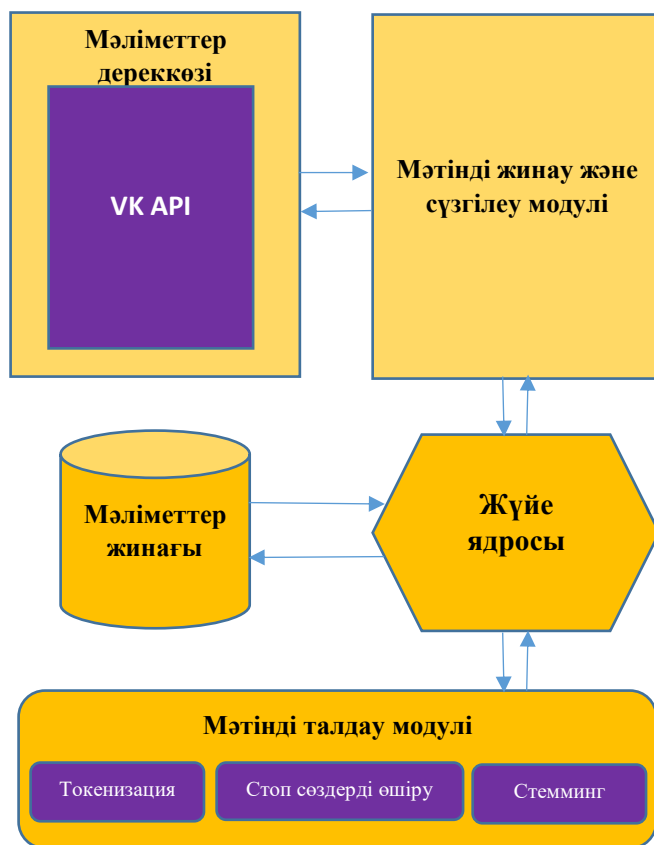
Бұл бөлімде терең оқыту алгоритмдері көмегімен табиғи тілді өңдеу әдістері негізінде әлеуметтік желідегі қазақ тіліндегі экстремистік бағыттағы мәтіндерді анықтауға арналған жүйенің жалпы архитектурасы келтіріледі. Машиналық және терең оқыту алгоритмдерін оқыту барысында қолданылған белгілерге сипаттама беріледі. Ұсынылған архитектура бірнеше модульдан тұрады (сурет 3.2): мәліметтерді алу және преппроцессинг модулі, белгілерді шығару және мәліметтерді белгілеу модулі және терең оқыту үлгілері арқылы талдау жүргізу модулі [113, 114].



Сурет 3.2 – Қазақ тіліндегі экстремистік мәтіндерді анықтау үшін ұсынылатын әдіс архитектурасы

### 3.1 Мәліметтерді алу және преппроцессинг модулі

Мәліметтерді алу үшін «ВКонтакте» әлеуметтік желісінің мәліметтер қорынан серверге https-сұраныстар арқылы қажетті ақпаратты алуға мүмкіндік беретін API Вконтакте интерфейсі қолданылды (сурет 3.3) [115].



Сурет 3.3 – Мәліметтерді алу және препроцессинг модулі

Әлеуметтік желі API-мен сұраныстар кітапханасы арқылы байланыс орнатылды. Жұмыс ортасы ретінде Pycharm Community Edition 2018 бағдарламалық жабдықтамасы қолданылды.

Мәтіндегі экстремистік бағытты анықтағанға дейін өңделмеген мәтін бірнеше қадамнан тұратын жіктеуге дейінгі дайындық кезеңдерінен өтуі керек. Берілген архитектура бойынша «ВКонтакте» әлеуметтік желісінен шикі хабарламаларды алу, алдын-ала өңдеу, тазалау, оқыту және тестілеу жинақтарын бөлу, оқыту алгоритмі мен мәтіндерді талдауға арналған табиғи тілді өңдеу үлгісін тестілеу, мәліметтер арасындағы шуды өшіру және мәтінді аннотациялау/белгілеу тапсырмалары орындалады. Мәліметтерді алдын ала өңдеу табиғи тілді өңдеу моделін құрастырудағы өте маңызды кезең болып табылады және бұл кезең үлгінің дәлдігіне тікелей әсер етеді. Бұл шикі мәтіндегі эмоджиларды, пунктуациялық белгілерді, http сілтемелер немесе басқа да әріптік-сандық белгілерге жатпайтын шулы символдарды өшіруді қамтиды. Егер біз кез келген мәтіндік корпусты қарастыратын болсақ, онда ең жиі кездесетін сөздердің ешбір тоналық күші жоқ стоп сөздер екендігін аңғарамыз. Сол себепті машиналық оқытудың ықтималдық үлгілерімен жұмыс істейтін болғандықтан, аталған сөздерді мәліметтер қорынан өшіру аса маңызды болып табылады. Қазақ тілінен өзге тілген жазылған барлық сөздер, сөйлемдер, сілтемелер, URL мекенжайлары және арнайы белгілер өшірілді.

Бұл жұмыста «Вконтакте» әлеуметтік желісінен жиналған мәліметтер пайдаланылды. Қосымша ақпаратты пайдаланбай дербес деректердің нақты дербес деректер субъектісіне тиесілігін анықтау мүмкін болмауы үшін дербес деректерді иесіздендірілді.

Мәліметтерге белгі тағайындау үшін эксперттердің көмегі пайдаланылды. Олар әр белгіні тағайындамас бұрын мәліметтер қорындағы әр сөйлемді мұқият оқыды. Бастапқыда мәліметтер қорындағы мәтіндер саны 30000 шамасында болды, алайда кей мәтіндер жіктеу сатысы үшін маңызды мәліметтерді қамтымағандықтан сарапшылар 27000-ға жуық мәтінді өшірді. Мәліметтер қорының көлемі 30000 мәтіннен 2600 мәтінге қысқарды. Ол жердегі мәтіндерге 1 (экстремистік) және 0 (бейтарап) белгілері тағайындалған.

Аталған белгілер эксперт мамандардың көмегімен белгіленді және келесі алгоритм бойынша тағайындалды:

- 1) Әр мәтінді оқып шығып, санат белгісін тағайындау;
- 2) Көп мағыналы мәтіндерді алып тастау;
- 3) Сәйкес мәтіндерді бір санатқа біріктіру, яғни жарылыс, соғысқа қатысу, өлтіру және т.с.с.
- 4) Әр мәтінге оқиғалардың екі түрінің бірін, яғни діни экстремизм және бейтарап санаттардың бірін тағайындау.

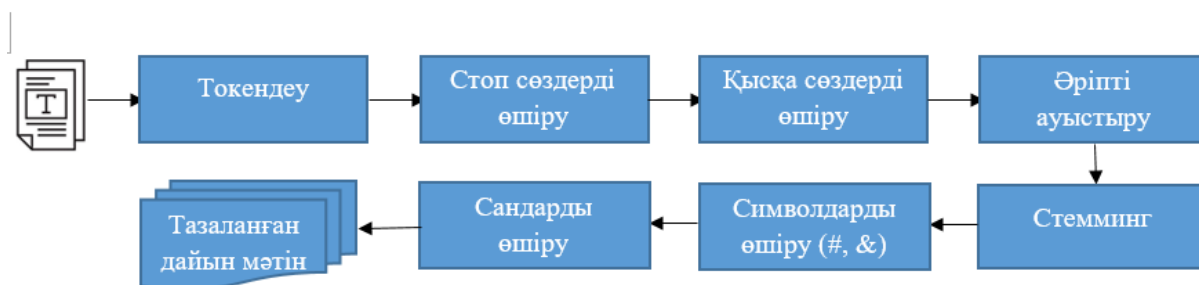
Әдетте деректердің екі түрі бар: құрылымдалған және құрылымданбаған. Құрылымдық деректерде оның мазмұнын, мағынасын немесе қолданылуын көрсететін әр деректер элементіне қатысты кілттер (атрибуттар, функциялар) бар. Құрылымдық деректердің типтік мысалы – дерекқордағы реляциялық кесте. Атрибут (баған) атауын және оның мәнін ескере отырып, осы мәнді қамтитын түйіндер (жолдар) жиынтығын жасауға болады.

Адамдар күнделікті өмірде пайдаланатын ақпарат негізінен құрылымданбаған түрде болады. Оның көп бөлігі табиғи тілде жазылған мәтіндік құжаттардан (кітаптар, журналдар, газеттер) тұрады.

Мәтінді алу үшін оны деректерді іздеу алгоритмдері қолдана алатын түрге келтіру керек. Олардың қатарына келесі препроцессинг қадамдарын жатқызуға болады: құжатты стандарттау, токенизация, стоп-сөздерді өшіру, стемминг немесе лемматизация және векторлық кеңістік кестесі [116]. Классикалық жіктеу және оқыту алгоритмдері мәтіндік құжаттарды бастапқы түрінде тікелей өңдей алмайды. Осылайша, алдын-ала өңдеу кезеңінде құжаттар басқарылатын түрге келтіріледі.

Деректер жиынтығында жартылай құрылымдалған/құрылымданбаған, көп мөлшердегі қажетсіз мәліметтер бар, олар болжам жасауда маңызды рөл атқармайды. Сонымен қатар, үлкен мәліметтер жиынтығы ұзақ уақыт оқытуды қажет етеді, ал стоп сөздер болжамның дәлдігін азайтады. Сондықтан мәтінді алдын-ала өңдеу есептеу ресурстарын үнемдеу үшін де, болжау дәлдігін арттыру үшін де қажет. Мәтінді алдын-ала өңдеу дәлірек болжауда маңызды рөл атқарады және модельдің жұмысын жоғарылатады. Сурет 3.4-те көрсетілгендей алдын-ала өңдеу кезеңінде жүзеге асырылады.





Сурет 3.4 – Препроцессинг қадамы

Жіктеу процессіне мағыналық жағынан ықпал етпейтін барлық сөздер стоп сөздер ретінде алынып тасталды, мысалы «және», «ертең», «кеше», «екен» және т.б. Деректерді тазалағаннан кейін және стоп сөздерді өшіргеннен кейін, әр сөйлем бос кеңістік негізінде сөздерге токенизацияланды.

Токенизация: бұл үздіксіз мәтінді сөздерге, символдарға және элементтерге бөлуді (токендер деп аталатын) қамтиды. Бұл кейінгі талдаудың орындалуына айтарлықтай әсер етеді, сондықтан бұл кезең дұрыс және тиімді орындалуы керек.

Стоп сөздерді өшіру: Келесі қадамда мәтіндегі стоп сөздер алынып тасталады. Стоп сөздер сөйлемдерді оқуды ыңғайлы етсе де, мәтінді талдауда қосымша маңыздылық бермейді. Стоп сөздерді алып тастау жіктеу алгоритмінің тиімділігін арттырады. Көп жағдайда оларды мәтіннен алып тастауға болады, өйткені олар кейінгі талдау үшін құнды ақпарат бермейді.

Қысқа сөзді өшіру: Мәтіндердегі ұзындығы үштен аспайтын қысқа сөздер алынып тасталады. Кей зерттеулерде оқыту алгоритмдерінің қысқа сөздермен дұрыс жұмыс істемейтіні және егер кіріс мәтінде қысқа сөздер болса, онда ол оның дәлдігіне әсер ететіндігін анықтады. Демек, жіктеуіштердің беріктігі мен тиімділігін арттыру үшін қысқа сөздер алынып тасталады.

Әріпті түрлендіру: қысқа сөздер өшірілгеннен кейін, әлеуметтік желі мәтіндеріндегі әріптер кіші әріпке айналдырылады. Бұл маңызды қадам, өйткені талдау алгоритмі регистрді ескереді. Ықтималдық модельдері, мысалы, «Жаман» және «жаман» сөздерін әр түрлі сөздер деп санайды және олар әр сөздің пайда болуын бөлек-бөлек санайды. Егер сөздер кіші әріпке айналдырылмаса, бұл жіктеуіштің тиімділігін төмендетуі мүмкін. Төменгі регистр – бұл мәтінді алдын-ала өңдеудің кең таралған әдісі. Төменгі регистрды орнату sklearn TFIDFVectorizer және Keras Tokenizer сияқты көптеген заманауи векторизаторлар мен токенизаторларда жасалады.

Стемминг: Стемминг – бұл аффикстерді сөздерден алып тастау және сөздерді олардың түбір формасына келтіру процесі. Мысалы, «соғысқа», «соғыста», «соғыстың» дегеніміз бірдей мағынадағы «соғыс» сөзінің морфологиялық өзгерістері болып табылады. Қосымшаларды алып тастау белгілердің күрделілігін азайтуға көмектеседі және жіктеуіштердің оқу қабілетін жақсартады.

@ және басқа да тыныс белгілерін өшіру: мәтінді алдын-ала өңдеудің тағы бір кең таралған әдісі – мәтіндік деректерден тыныс белгілерін жою. Бұл тағы да мәтінді стандарттау процесі, ол "ура" және «ура!» тіркестерін бір тіркес ретінде тану үшін қажет. Сондай-ақ, пайдалану жағдайына байланысты алып тастау үшін тыныс

белгілерінің тізімін мұқият таңдау керек. Мысалы, python-да string.punctuation келесі тыныс белгілерін қамтиды: !"#\$%&\'()\*+,-./:;<=>?@[\\]^\_`{|}~`. Біздің қажеттіліктерімізге сәйкес көбірек тыныс белгілерін қосуға немесе жоюға болады.

Келесі қадамда мәтіндегі сандық мәндер алынып тасталады, өйткені олар мәтінді талдау үшін ешқандай мәнге ие емес, ал оларды жою модельдерді оқытудың күрделілігін төмендетеді.

### **3.2 Белгілерді шығару және мәліметтерді белгілеу модулі**

Мәтін – ақпарат беру үшін жиі қолданылатын медиа. Көлемді мәтіндердің пайда болуымен ақпараттың шамадан тыс жүктелуі және деректердің артық болуы сияқты проблемалар барған сайын жиі байқалуда. Көрнекі тұрғыдан алғанда, мәтінді визуализациялау – бұл ең маңызды кезең. Осылайша, мәтінді визуализациялау технологиясы мәтіндегі мазмұнды тез алу үшін адамдар визуалды қабылдауға тән параллельді өңдеу мүмкіндіктерін қолдана алатындай етіп, мәтінде айту қиын болатын мазмұн мен ережелерді көрнекі таңбалар түрінде білдіреді.

Мәтінді визуализациялау табиғи тілді өңдеуге негізделген, сондықтан мәтінді талдаудың кең таралған әдістері – сөз теру моделі, аталған нысандарды тану, кілт сөздерді шығару, тақырыпты талдау, тоналдылықты талдау және т. б. Мәтінді талдау процесі негізінен сөздерді сегментациялау, шығару және қалыпқа келтіру сияқты әрекеттерді қолдана отырып, сөздік деңгей мазмұнын шығаратын функцияларды алуды қамтиды, сонымен қатар векторлық кеңістік моделін құру үшін функцияларды қолданады және оны төмен өлшемді кеңістікте көрсету немесе тақырыптарды пайдалану үшін өлшемді азайтады. Модель сипаттамаларды өңдейді және соңында өңделген деректерді визуалды бейнелеу үшін икемді және тиімді түрде ұсынады.

Белгілерді құрастыру – машинаны оқытуға қажетті белгілерді құру үшін пәндік сала мәліметтерін қолдану процесі. Белгілерді құрастыру машиналық оқыту қосымшаларының негізі болып табылады және ол қиын, әрі көп еңбекті талап ететін шара.

Белгі – біз талдайтын немесе болжам жасайтын барлық тәуелсіз нысандарға тән қасиет. Үлгі үшін қажетті кез келген белгі пайдалы болуы мүмкін. Белгі есепті шешуге көмектесетін сипаттама болып табылады.

Мәліметтердегі белгілер болжам жасайтын үлгілер үшін маңызды болып табылады және келешекте алынатын нәтижеге әсер етеді. Белгілердің саны мен сапасы үлгінің сапасына, оның дұрыс болжам жасауына үлкен әсерін тигізеді. Белгі жақсы болған сайын нәтиже де соғұрлым дәлірек болады деуге болады, алайда кей жағдайда нәтиже тек таңдалған белгілерге ғана емес, сонымен қатар үлгі мен мәліметтерге де байланысты болады. Дегенмен, дұрыс белгілерді таңдау өте маңызды.

Белгілерді оқыту – жүйеге белгілерді анықтауға немесе бастапқы (шикі) мәліметтерді жіктеуге қажетті бейнелеулерді автоматты түрде анықтауға мүмкіндік беретін техникалар жинағы. Бұл үдеріс белгілерді қолмен құрастыру есебін алмастырады және машинаға белгілерді зерттеуге де, белгілі бір есептерді шешу үшін оларды пайдалануға да мүмкіндік береді.

Белгілерді оқыту жіктеу сияқты машиналық оқыту есептерінде математикалық және есептеу түрінде өңдеу ыңғайлы болатын кірістің қажет болуымен түсіндіріледі.

Алайда сурет, бейнежазба және тетіктердің жазбалары сияқты шынайы мәліметтер аталған есептердің алгоритмдік анықтамасына сай келмейді.

NLP-дағы ең маңызды процесстердің бірі белгілерді таңдау (feature selection) болып табылады. Белгілерді таңдау процесінде модельдің тиімділігі мен өнімділігін арттыру үшін корпустан өзекті және ең пайдалы белгілер алынады. Мәтіндік деректерді бастапқы форматта математикалық алгоритмдерге қолдану мүмкін болмағандықтан, оларды математикалық түрде векторлық форматта ұсынуымыз керек. Белгі мәтінді бірқатар ерекшеліктерге ие векторлық кеңістікке айналдырады. Мәтін форматындағы мәліметтерде көп шу қамтылғандықтан, олар категориялық мәліметтер болып табылады және оны машиналық оқыту негізінде тікелей қолдану мүмкін емес, сондықтан математикалық модельді оқыту кезінде бұл кірістерді тиімді өңдей алатындай етіп мәтіндік деректерді сандық векторлық форматқа ауыстыруымыз керек. Шикі мәтіндік деректерді белгі векторларына (feature vectors) айналдыру белгіні ұсыну (feature representation) деп аталады.

Мәтіндік деректерді векторлық форматта ұсыну үшін сөздердің пакеті Bag of Words (BOW), TF-IDF сияқты көптеген тәсілдер бар. Бұл векторлар тек лексиконға негізделген векторлар, әдетте құжаттағы сөздердің санын немесе оның салыстырмалы салмағын білдіреді. Сөздерді векторлық бейнелеу (word embedding) – тілді модельдеу және табиғи тілді өңдеудегі қандай да бір сөздіктегі сөздерге кішігірім өлшемдегі векторларды сәйкестендіруге бағытталған бейнелеулерді оқытуға арналған әр түрлі тәсілдер кешені.

Кіріс корпусы – пайдаланушылардың хабарламасын білдіретін  $D = \{d_1, d_2, d_3 \dots \dots, d_n\}$  құжаттар жиынтығы және әр  $d_i$  құжатта қолданылатын сөздер саны бойынша әр түрлі ұзындықта болуы мүмкін. Белгілер матрицасында бұл құжаттар ұзындығы  $m = |V|$  жоғары өлшемді векторлар ретінде ұсынылған, мұндағы  $m$  – сөздік қорының мөлшері.

$d_i = \{w_1, w_2, w_3, \dots \dots w_m\}$ , мұндағы  $w_i$   $d_i$  құжатында  $w_i$  сөзінің бар-жоқтығын көрсететін екілік мәнге ие бола алады.

Белгі векторлары – бұл мәтіннің сандық көрінісі. Бұл машиналық оқыту классификаторы өңдей алатын енгізудің нақты формасы, мәтінді өңдеу үшін қолданылатын белгілерді қалыптастырудың бірнеше әдістері бар. Векторларды жасаудың келесі ерекшеліктері қолданылды.

### 3.2.1 TF-IDF

Мәтіндік деректерді сандарға ауыстыру қажет және мәтіндік деректерді сандарға өңдеудің кең қолданылатын әдісі – TF-IDF болып табылады. TF-IDF-де мәтіндік мәліметтер векторларға айналады, олар сөздердің реттілігінің нақты дәйектілігін ескермейді. Корпустағы әрбір сөз TF-IDF санымен байланысты, бұл әр сөздің корпус үшін қаншалықты маңызды екендігін көрсетеді.

Сөздерді сандарға айналдырғаннан кейін, TF-IDF сандық мәндері машиналық оқыту әдістері түсіндіре алатын контекстте бақыланатын оқыту жіктеуіштеріне беріледі.

TF-IDF – құжаттағы терминдердің өлшемді векторлық көрінісі. Ол құжатта пайда болған әр терминге салыстырмалы салмақты терминнің дифференциалдау

қабілеттілігіне сәйкес тағайындайды, бұл құжатқа тиісті белгіні тағайындауға көмектеседі. TF-IDF бойынша, егер термин корпустағы барлық құжаттарда пайда болса, онда ол онша маңызды емес және оған аз салмақ беру керек, екінші жағынан салыстырмалы түрде аз құжаттарда пайда болса, оған үлкен салмақ беру керек деп есептеледі.

TF (term frequency — сөз жиілігі) — қандай да бір сөздің құжатта кездесуінің құжаттағы барлық сөздер санына қатынасы.  $t_i$  сөзінің жеке құжат шеңберіндегі маңыздылығы есептеледі.

$$tf(t, d) = \frac{n_t}{\sum_k n_k} \quad (3.1)$$

Мұндағы  $n_t$  –  $t$  сөзінің құжатта кездесу саны, ал бөлімінде берілген құжаттағы сөздердің жалпы саны.

IDF (inverse document frequency — құжаттың кері жиілігі) — қандай да бір сөздің жинақ құжаттарында кездесу жиілігінің инверсиясы. IDF кең қолданылатын сөздердің салмағын азайтады. Нақты бір құжаттар жинағындағы әрбір бірегей сөз үшін бір ғана IDF мәні болады.

$$idf(t, D) = \log \frac{|D|}{|\{d_i \in D \mid t \in d_i\}|} \quad (3.2)$$

Мұндағы  $|D|$  - жинақтағы құжаттар саны,  $|\{d_i \in D \mid t \in d_i\}|$  -  $D$  жинағындағы  $t$  кездесетін құжаттар саны ( $n_t \neq 0$ ).

TF-IDF өлшемі екі шаманың көбейтіндісі түрінде есептеледі [117]:

$$tf - idf(t, d, D) = tf(t, d) * idf(t, D) \quad (3.3)$$

### 3.2.2 n-грамдар

Берілген жұмыста қолданылған тағы бір белгі – n-грамдық ендіру әдісі. n-gram тиімді техника болып табылады, өйткені ол корпустағы сөздер ретін ұсынады.

N-gram – символдар немесе сөздер қатарынан тұратын токен. Токен мәтіндер корпусы бойымен сырғымалы терезені жылжыту арқылы қалыптасады, терезенің өлшемі токеннің өлшеміне байланысты, ал жылжыту кезең бойынша орындалады, әр кезең сөзге немесе символға сәйкес келеді.

[118] жұмыста n-грам бірнеше сөздің тізбегі ретінде анықталған: 1 грам (униграмма), 2 грам немесе биграмма – екі сөзден құралған тізбек, ал 3 сөз немесе триграмма – үш сөзден құралған тізбек.

Мәтіндік деректерді салмақталған векторға айналдыру және хабарламадағы сөз тізбегіне ықтималдықтарды тиімді бөлу үшін n-граммдық модель қолданылады. N-грамм моделін түсіну үшін хабарламаны мысал ретінде келтірейік, мысалы: «Менің әпкем ауырып жатыр, ол тезірек сауығып кетеді деп үміттенемін!». Биграммдық n-граммдық интерпретациясы (бұл жағдай үшін  $N - 1 = 1$ , сөздің пайда болуын алдыңғы сөзге сүйене отырып болжайды) «менің әпкем», «әпкем ауырып», «ауырып жатыр», «жатыр ол», «ол тезірек», «тезірек сауығып», «сауығып кетеді», «кетеді деп», «деп үміттенемін» түрінде болады.

$d_1$  және  $d_2$  құжаттарының ұқсастығын олардың  $n$ -gram түрінде ұсынулары  $S_n(d_1)$  және  $S_n(d_2)$  арқылы анықтауға арналған Жаккард коэффициенті келесі теңдеуде келтірілген:

$$\text{sim}(d_1, d_2) = \frac{|S_n(d_1) \cap S_n(d_2)|}{|S_n(d_1) \cup S_n(d_2)|} \quad (3.4)$$

$d_1$  және  $d_2$  құжаттары бір-бірінің көшірмесі екендігін анықтау үшін шектік мән пайдаланылады. Әрбір қосымша үшін терезенің өлшемі мен ұқсастық шегі тәжірибелер барысында таңдалады [119].

[120] жұмыста  $n$ -грам негізіндегі модельді  $n$ -грамнің тілдегі ықтималдығын болжау үшін пайдаланылатын ықтималдық моделі деп атаған. Берілген ықтималдықтар негізінде  $d$  құжатының белгілі бір  $c$  классына тиіс екендігін анықтауға болады:

$$P(c|d) = \prod_{i=1}^n P(w_i | w_{i-n+1}^{i-1}) \quad (3.5)$$

$$P(w_i | w_{i-n+1}^{i-1}) = \frac{\text{freq}(w_{i-n+1}^{i-1} w_i)}{\text{freq}(w_{i-n+1}^{i-1})} \quad (3.6)$$

мұндағы  $\text{freq}$  -  $n$ -gram-ның класстағы жиілігін білдіреді.

### 3.2.3 Bag-of-words

Құжаттар оқыту процесіне жарамды болатын түрде ұсынылуы керек. Мәтінді жіктеуде қолданылатын әдеттегі мәтінді ұсыну схемасы – бұл bag-of-words, ол жерде синтаксис, сөздердің рейтингі және мәтіннің грамматикалық ережелері ескерілмейді. Бұл схема мәтіндік құжаттардың компоненттері арасындағы мағыналық қатынастарды құра алмайды. Берілген әдісті пайдалану барысында құжат әрқайсысы белгілі бір мәліметтер жиынтығынан сөздіктегі терминге немесе сөз тіркесіне сәйкес келетін белгілер векторы түрінде ұсынылады. Белгінің әрбір элементінің мәні терминнің белгінің нақты өлшеміне сәйкес құжаттағы маңыздылығын көрсетеді.

Термин – семантикасы құжаттың негізгі тақырыптарын есте сақтауға мүмкіндік беретін сөз. Әр терминнің салмағы болады.  $D$  құжаттар жинағы үшін  $V = \{t_1, t_2, \dots, t_{|V|}\}$  – жинақтағы бірегей терминдер жиынтығы, ал  $t_i$  – термин болады.  $V$  жинағы әдетте жинақ сөздігі деп аталады,  $|V|$  - оның көлемі, яғни  $V$ -дағы терминдер саны.  $d_j \in D$  құжатының әрбір  $t_i$  терминінің  $w_{ij} > 0$  салмағы болады,  $d_j$  құжатында кездеспейтін терминнің салмағы 0-ге тең,  $w_{ij} = 0$ . Осылайша, әрбір  $d_j$  құжаты  $d_j = (w_{1j}, w_{2j}, \dots, w_{|V|j})$  терминдер векторы арқылы ұсынылады, мұндағы әрбір  $w_{ij}$  салмағы  $t_i \in V$  терминіне сәйкес келеді және  $t_i$  терминнің  $d_j$  құжаттағы маңыздылығын анықтайды [119, б.215].

Bag-of-words белгілердің ортогоналдылығына болжам жасайды және сөздердің реті мен грамматиканы елемейді.  $C = \{c_1, c_2, \dots, c_{|C|}\}$  – қызығушылық тудыратын класстар, ал  $F = \{f_1, f_2, \dots, f_{|F|}\}$  – мәліметтер жинағындағы кем дегенде бір құжатта кем дегенде бір рет кездесетін бірегей функциялар жиынтығы,  $T_r = \{d_1, d_2, \dots, d_{|T_r|}\}$  – оқытуға арналған мәтіндік құжаттар жинағы, ал  $w(f_i, d_j) = d_j$  құжаттағы  $f_i$  нысанының салмағы болсын.  $d_j$  құжаты нысанның  $d_j = [w(f_1, d_j), \dots, w(f_{|F|}, d_j)]$  салмақ векторы арқылы ұсынылады. Берілген  $w(f_i, d_j)$  салмағы  $d_j$  құжатын семантикалық сипаттау кезіндегі  $f_i$  белгісінің маңыздылығын анықтайды [121].

Мәтіндік құжаттар арасындағы ұқсастық bag-of-words ұсыну негізінде анықталады және мәліметтер қорындағы әрбір құжат пен қолданушы сұранысы көпөлшемді вектор түрінде көрсетілетін векторлық кеңістік моделін пайдаланады.

BOW (Bag of Word) – мәтіндік мәліметтерді N-gram арқылы құжаттар матрицасына түрлендіретін үлгі. N-грамм ықтималдықты тағайындау үшін қарастырылатын сөздердің санына байланысты униграмма, биграмм немесе триграмма болуы мүмкін.

BOW негізіндегі ықтималдық модельдер көбінесе униграмма ерекшеліктерімен жоғары дәлдікке жетеді, бірақ әдетте сөздерді олардың арасындағы қатынастарды ескермей оқшаулайды. BOW-да контекстік ақпаратты қарастырғанмен, биграмм мен триграмма қолданылады, бірақ сөздердің ұзақ мерзімді тәуелділігі жоғалады.

### 3.2.4 Word2vec

Word2vec – сөз векторларының қысылған кеңістігін құрудың нейрондық желілерді пайдаланатын тәсілі. Ол кіріс ретінде үлкен көлемді мәтіндік корпусты қабылдайды және әр сөзге векторды сәйкестендіреді. Берілген алгоритм бойынша алдымен сөздік құрылады, содан кейін сөздердің векторлық бейнелеуі есептеледі. Векторлық бейнелеу контекстік жақындыққа негізделеді: мәтінде бірдей сөздердің қасында кездесетін (яғни мағыналары ұқсас) сөздер векторлық бейнелеуде жоғары косинустық ұқсастыққа ие болады:

$$\text{similarity}(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}} \quad (3.7)$$

Жақында жүргізілген зерттеулерде мәтіндік құжаттарды нейрондық тілдік модельдерді қолдана отырып өңдеу табиғи тілді өңдеу (NLP) тапсырмаларын, оның ішінде сентимент талдау жасауды, лингвистикалық модельдеуді және машиналық аударманы болжау нәтижелерін едәуір арттыра алатынын көрсетеді. Соңғы уақытта мәтіндік бейнелеуді үйрену үшін нейрондық тілдік модельдер қолданылуда, бұл үлкен құрылымданбаған мәтіндік мәліметтерден алынған сөздерді ендіру (word embedding) схемаларының тиімділігіне байланысты болып табылады.

Сөздерді ендіру (word embedding) схемалары аз өлшемді мәтіндік құжаттарды тиімді түрде ұсынуға мүмкіндік береді, осылайша bag-of-words схемасында кездесетін сирек және үлкен өлшемділік проблемаларын жояды, олар мәтіндік мәліметтердің тек синтаксистік қана емес, сонымен қатар семантикалық та ақпаратын қамтиды. Мәтіндік құжаттарды терең оқыту алгоритмдері негізінде жіктеу әдістері белгілерді қалыптастыру кезеңі автоматты түрде орындалатындықтан дәстүрлі машиналық оқыту алгоритмдерімен салыстырғанда тиімді болып келеді.

Word2vec моделі – бұл ендіру қабатынан, шығыс қабатынан және жасырын қабаттан тұратын жасанды нейрондық желіге негізделген сөз ендіру схемасы. Ол белгілі бір сөздің басқа сөздермен еніп кету ықтималдығын анықтау арқылы сөздерді ендіруді үйренуге бағытталған. Модельдің екі негізгі архитектурасы бар, skip-gram (SG) және continuous-bag-of words (CBOW). CBOW архитектурасы мақсатты сөзді әр сөздің мазмұнын кіріс ретінде қабылдау арқылы анықтайды; SG архитектурасы, керісінше, мақсатты сөзді кіріс ретінде қабылдау арқылы мақсатты сөздің айналасындағы сөздерді болжайды. CBOW архитектурасы мәліметтердің аз мөлшерімен дұрыс жұмыс істей алады.

Осы техникалар арқылы пайда болған векторлар нейрондық желілердің кірістірілген қабатына кіріс ретінде беріледі. Кірістірілген қабаттардың нәтижесі

терең оқыту модельдерінің келесі толық қосылған қабатына беріледі. Тестілеу / тексеру кезеңінде модельді өңдеу аяқталғаннан кейін әр мәтінге екі класстан тиісті класс белгісі тағайындалады.

### 3.3 Терең оқыту үлгілері арқылы талдау жүргізу модулі

#### 3.3.1 Терең оқыту үлгілері

Терең оқыту (ағылш. Deep learning) – машиналық оқыту әдістерінің жиынтығы (оқытушымен, оқытушыны ішінара тарта отырып, оқытушысыз, күшейтілген). Терең оқытудың көптеген әдістері жеткілікті өнімділікке ие және бұрын тиімді шешуге мүмкіндік бермейтін көптеген мәселелерді шешуге мүмкіндік береді, мысалы, компьютерлік көру, машиналық аударма, сөйлеуді тану есептерінде жиі пайдаланылады.

Терең нейрондық желі (ағылш. DNN – deep neural network) – кіру және шығу қабаттары арасындағы бірнеше қабаты бар жасанды нейрондық желі. Терең нейрондық желі сызықтық немесе сызықтық емес корреляцияға қарамастан кірістерді шығысқа айналдыру үшін математикалық түрлендірудің дұрыс әдісін табады. Желі қабаттар бойынша қозғалады, әр қабат үшін шығу ықтималдығын есептейді. Пайдаланушы нәтижелерді қарап, желіні көрсететін ықтималдылықты таңдай алады (мысалы, белгілі бір шектен жоғары) және ұсынылған белгіні желіге қайтара алады. Әрбір математикалық түрлендіру қабаты болып саналады, ал күрделі терең нейрондық желіде көптеген қабаттар бар, сондықтан "терең" желілер деп аталады.

Терең оқыту – бұл көптеген сызықты емес түрлендірулерді қолдана отырып, жоғары деңгейлі абстракцияларды модельдеуге арналған машиналық оқыту алгоритмдері. Ең алдымен терең оқытуға келесі әдістер мен олардың өзгерістері жатады [122-126]:

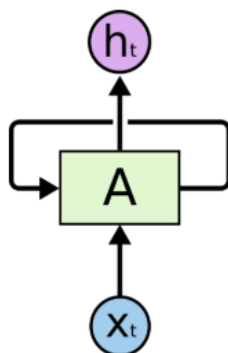
- 1) оқытушысыз оқытудың белгілі бір жүйелері, мысалы, Больцманның алдын-ала оқыту машинасы, авто-кодтаушы, терең сенім желісі, генеративті-қарсыласу желісі,
- 2) оқытушымен бірге оқытудың белгілі бір жүйелері, мысалы, үлгіні тану технологиясын жаңа деңгейге көтерген конвульсиялық нейрондық желі,
- 3) уақыт өте келе процестерде оқуға мүмкіндік беретін қайталанатын нейрондық желілер,
- 4) тізбек элементтері мен тізбектер арасындағы кері байланысты қосуға мүмкіндік беретін рекурсивті нейрондық желілер.

Осы әдістерді біріктіре отырып, жасанды интеллекттің әртүрлі міндеттеріне сәйкес келетін күрделі жүйелер жасалады

#### 3.3.2 Рекуррентті нейрондық желілер

Рекуррентті нейрондық желілер (ағылш. Recurrent neural network; RNN) – элементтер арасындағы байланыс бағытталған тізбекті құрайтын нейрондық желілердің түрі. Оның көмегімен уақыт өте келе оқиғалар сериясын немесе дәйекті кеңістіктік тізбектерді өңдеуге болады.

Көпқабатты перцептрондардан ерекшелік ретінде рекуррентті желілер кез келген ұзындықтағы тізбектерді өңдеу үшін өзінің ішкі жадын пайдалана алады. Сол себепті RNN қандай да бір тұтас нәрсе бірнеше сегментке бөлінетін, яғни қолмен жазылған мәтінді тану, тілді тану сияқты есептерді шешуде жиі пайдаланылады. Рекуррентті желілердің бірнеше архитектуралық шешімдері бар. Соңғы уақытта ұзақ және қысқа мерзімдік жады (LSTM) және басқарылатын рекуррентті блок желісі танымалдылыққа ие (сурет 3.5) [127].

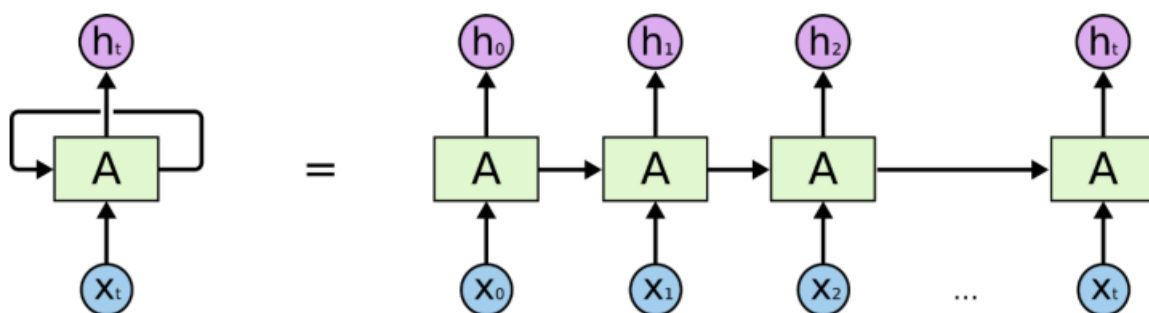


Сурет 3.5 – Қарапайым рекурренттік желі

Суреттегі нейрондық желінің А бөлігі қандай да Х мәліметтерін кіріс ретінде қабылдайды және қандай да Н мәнін шығыс ретінде қайтарады. Циклдық байланыс желінің алдыңғы қадамындағы ақпаратты келесі қадамға беруге мүмкіндік береді.

Рекуррентті нейрондық желілердің сан алуан түрлері, шешімдері және құрылымдық элементтері бар. Рекурренттік желінің қиындығы егер әрбір уақыт қадамын есепке алу қажет болса, онда әрбір уақыт қадамы үшін жеке нейрондар қабатын құру қажеттілігінде болып табылады, ал бұл өте үлкен есептеу күрделілігін туындатады. Егер есептеуді тіркелген уақыт терезесімен шектейтін болсақ, онда алынған үлгілер трендтің ұзақ мерзімділігін көрсетпейді.

Рекуррентті нейрондық желілер қарапайым нейрондық желілерден аса ерекшеленбейді. Оларды бір желінің бірнеше көшірмесе ретінде елестетуге болады, әрбір көшірме келесі көшірмеге хабарлама жібереді деп есептеледі. Егер циклды ашатын болсақ, ол келесідей көрініске ие болады (сурет 3.6):



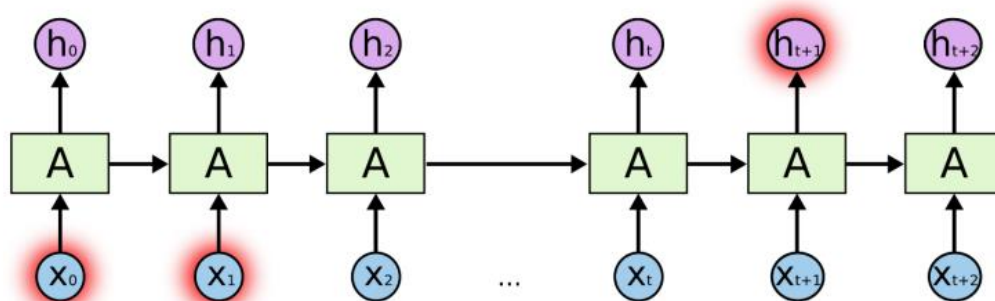
Сурет 3.6 – Рекурренттік желі циклы

Мұндай «тізбек» рекуррентті нейрондық желілердің тізбектер мен тізімдермен тығыз байланысты екендігіні көрсетеді. Желі әрқайсысы басқа түйіндермен



байланысатын түйіндерден тұрады. Уақыт өте келе әр нейронның активация шегі өзгеріп отырады және ол нақты сан болып табылады. Әрбір байланыстың айнымалы нақты салмағы бар. Түйіндер кіріс, шығыс және жасырын болып бөлінеді.

Рекуррентті нейрондық желілердің негізгі артықшылықтарының бірі – ағымдағы есеп үшін бұрын алынған ақпаратты пайдалану мүмкіндігі. Бірақ кей жағдайда ақпарат пен ол қажет етілетін орын арасындағы ара қашықтық алыс болуы мүмкін, өкінішке орай, ара қашықтық ұзарған сайын рекуррентті нейрондық желілер ақпаратты байланыстыруды оқытуға қабілетсіз бола бастайды (сурет 3.7).

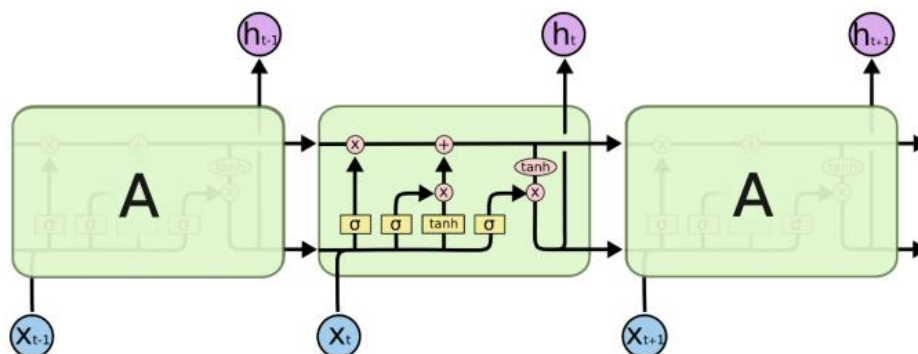


Сурет 3.7 – Рекурренттік желілерде кездесетін мәселе мысалы

Теориялық тұрғыдан алғанда рекуррентті нейрондық желілер мұндай ұзақ мерзімді тәуелділіктерді өңдей алады, алайда тәжірибелік тұрғыда олар мұндай есептер үшін оқытыла алмайды[127, б.78].

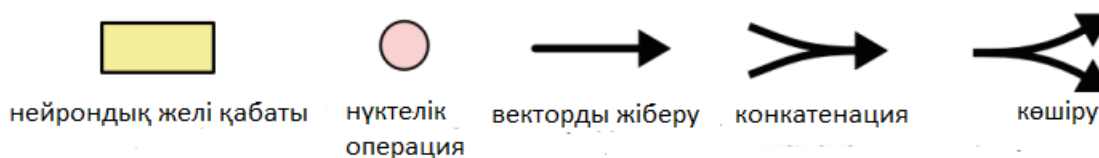
### 3.4 LSTM желілері

Ұзақ-қысқа мерзімді жады (Long Short Term Memory, LSTM) – рекуррентті нейрондық желілердің ұзақ мерзімді тәуелділіктерді оқуға қабілетті ерекше түрі. Олар көптеген есептерді шешу үшін тиімді пайдаланылады. LSTM ұзақ мерзімді тәуелділіктер мәселесінің алдын алу үшін арнайы жобаланған. LSTM желісі «ұмыту» ұяшығы деп аталатын рекуррентті ұяшықтар арқылы басқарылады. Қателіктер шектелмеген виртуалды қабаттар арқылы уақыт бойынша кері қайтарылады. LSTM-дегі оқыту осылай орындалады, сонымен қатар мыңдаған өткен уақыт мезеттері жайлы ақпарат жадыда сақталады. LSTM желісі 4 нейрондық қабаттан тұрады және олар өзара байланысты болып табылады (сурет 3.8) [127, б.79].



Сурет 3.8 – LSTM желісінің архитектурасы

Мұнда келесідей белгілеулер енгізілген (сурет 3.9):

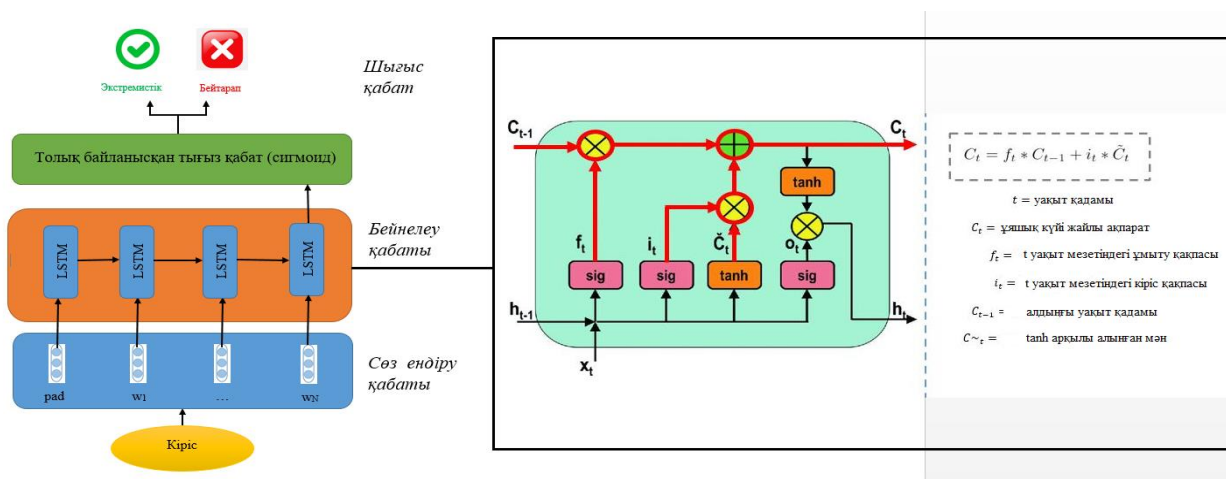


Сурет 3.9 – LSTM желісіндегі белгілеулер

Суреттегі әрбір сызық бір түйіннің шығысындағы векторды келесі түйіндерге кіріс түрінде жібереді. Қызғылт дөңгелектер векторларды қосу сияқты нүктелік операторлар, ал сары тіктөртбұрыштар нейрондық желінің оқытылған қабаттары болып табылады. Бірігетін сызықтар конкатенацияны, тармақталатын сызықтар олардың көшірілетіндігін және көшірменің әр түрлі орынға жіберілетінін білдіреді.

### 3.5 Ұсынылатын модель

Ұсынылатын модель – иерархиялық көп сатылы процесс және бірнеше қабаттан тұрады. Бірінші қабат – бұл алдын-ала дайындалған кірістіруді қолдана отырып, сөздерді векторлық кеңістікке бейнелейтін сөз ендіру қабаты. Әрі қарай бейнелеу қабаты орналасады, ол әр мәтінді бейнелеу матрицасын алу үшін LSTM-ді пайдаланады, содан кейін мәтіндер туралы онтайлы көрініс беру үшін контекст пен ағымдағы сөйлемді бейнелеуді біріктіретін әртүрлі әдістерді қолданады. Келесі қабат – мәтінді жіктеу нәтижесін беретін шығыс деңгей (сурет 3.10).



Сурет 3.10 – Ұсынылатын модель архитектурасы

#### 3.5.1 Word Embedding Layer – Сөз ендіру қабаты

Сөз ендіру қабаты әр сөзді мағыналық және синтаксистік ақпаратты жинақтай алатын үлкен векторлық кеңістікке бейнелеуге жауапты. Матрицаның әр бағанында сәйкес сөздің ендірілген сөзі сақталады.

Сөздерді векторизациялау (word embedding) – сөздер мен құжаттарды векторлық бейнелеу арқылы ұсынуға арналған әдістер санаты. Вектор сөздің үзіліссіз кеңістікке проекциясын білдіреді. Сөзді векторлық кеңістікте бейнелеу мәтіннен

алынады және сөзді пайдалану барысында оның маңында орналасатын сөздерге негізделеді.  $X_1. . . x_T$  кіріс сөйлемдегі сөздердің тізбегін білдірсін деп есептейік. Біріншіден, біз әр сөздің бекітілген сөзін ендіру үшін сөз векторларын қолданамыз. Осы қабат арқылы сөйлем матрица түрінде ұсынылады:  $X \in R_m \times T$ , мұндағы  $m$  – сөз векторының өлшемі, ал  $T$  – сөйлемнің ұзындығы. Біріншіден, сөздерді ендіру модельдерін қолдану арқылы хабарламалардағы сөздер векторларға түрлендіріледі. Содан кейін, сөздер арасындағы контекстік тәуелділіктің қашықтығын зерттеу үшін сөйлемдердегі сөз тізбектері LSTM-ге енгізіледі.

Сөз ендіру қабаты кездейсоқ сандармен инициалданады және оқыту жинағындағы барлық мәліметтер үшін векторизациялау тапсырмасын орындайды. Бұл қабат желідегі алғашқы жасырын қабат және оның үш аргументі бар:

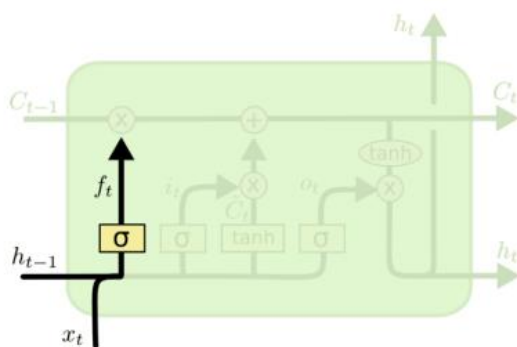
input\_dim: мәтіндік мәліметтер сөздігінің өлшемі.

output\_dim: сөздер векторизацияланатын векторлық кеңістіктің өлшемі. Ол әрбір сөз үшін осы қабаттағы шығыс векторларының өлшемін анықтайды.

input\_length: кіріс тізбектердің ұзындығы.

### 3.5.2 LSTM-Based Representation – Бейнелеу қабаты.

LSTM желісіндегі негізгі элемент – ұяшық күйі, яғни диаграмманың жоғары жағынан өтетін көлденең сызық. LSTM ұяшық күйіндегі ақпаратты өшіруге немесе оған жаңа ақпаратты қосуға қабілетті, алайда аталған қабілет тетік деп аталатын құрылым арқылы басқарылады. Олар сигмоидтық нейрондық желілерден және нүктелік көбейту операцияларынан тұрады. LSTM-дегі алғашқы қадам ұяшық күйіндегі қай ақпаратты алып тастау қажеттігі туралы шешім қабылдау болып табылады. Бұл шешім «ұмыту тетігі» деп аталатын сигмоидтық қабат арқылы қабылданады. Оның кірісіне  $h_{t-1}$  және  $x_t$  мәндері беріледі және шығысында 0 мен 1 арасындағы сан шығарылады (сурет 3.11).

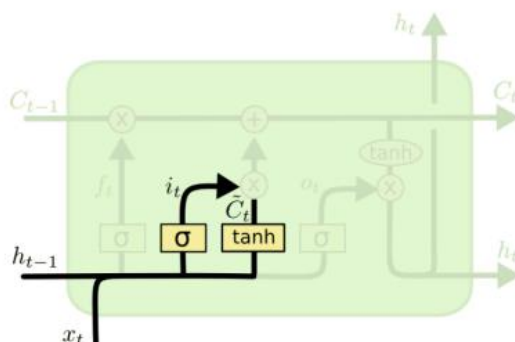


Сурет 3.11 – LSTM желісіндегі «ұмыту тетігі»

$$f_t = \sigma(W_f * [h_{t-1}, x_t] + b_f) \quad (3.8)$$

Келесі қадам – ұяшық күйінде қандай жаңа ақпаратты сақтау керектігі жайлы шешім қабылдау. Бұл қадам екі бөліктен тұрады. Біріншіден, «кіріс тетігі» деп аталатын сигмоидтық қабат қандай мәндердің жаңартылатындығы жайлы шешім қабылдайды. Әрі қарай, гиперболалық тангенс қабаты күйге қосылатын жаңа  $\sigma_t$

мәніне үміткерлер векторын құрады және күй үшін жаңартуды алу үшін аталған екі бөлік біріктіріледі (сурет 3.12).

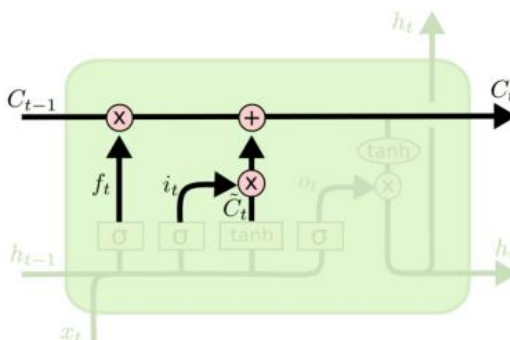


Сурет 3.12 – LSTM желісіндегі «кіріс тетігі»

$$i_t = \sigma(W_i * [h_{t-1}, x_t] + b_i) \quad (3.9)$$

$$\tilde{C}_t = \tanh(W_c * [h_{t-1}, x_t] + b_c) \quad (3.10)$$

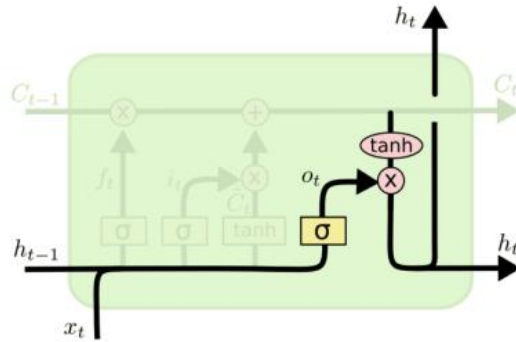
Әрі қарай ескі  $\sigma_{t-1}$  ұяшық күйін жаңа  $\sigma$  мәнімен жаңартатын уақыт келеді (сурет 3.13) .



Сурет 3.13 – LSTM желісіндегі «жаңарту тетігі»

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t \quad (3.11)$$

Келесі қадамда шығыс ретінде шығарылатын мәндер жайлы шешім қабылданады. Алдымен шығысқа ұяшық күйінің қандай бөліктері жіберілетіні жайлы шешім қабылдайтын сигмоидтық қабат іске қосылады. Әрі қарай ұяшық күйі гиперболалық тангенс қабатынан өтеді (мәндерді -1 және 1 аралығына сыйдыру үшін) және ол симгмоидтық тетік шығысына көбейтіледі (сурет 3.14) [127, б.80].



Сурет 3.12 STM желісіндегі «шығыс тетігі»

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o) \quad (3.12)$$

$$h_t = o_t + \tanh(C_t) \quad (3.13)$$

Берілген LSTM нейрондық желісіндегі гиперболалық тангенс активациялау функциясы келесі формула бойынша есептеледі:

$$\tanh(x) = \frac{2}{1+e^{-2x}} - 1 \quad (3.14)$$

мұндағы  $i, f, o, c$  – кіріс күйі, ұмыту күйі, шығыс күйі, жады ұяшығы,  $x$  – ұяшықтарды активациялаудың кіріс векторы,  $h$  – жасырын вектор,  $W_i$  – жасырын кіріс элементтер матрицасы,  $W_o$  – кіріс және шығыс элементтер матрицасы,  $b$  – базалық вектор.

Бұл архитектурада бастапқы сөйлемнің тереңдетілген бейнесі құрылады. LSTM қабаттарының ақырғы нәтижесі бір матрицаға біріктіріледі, бұл матрица толық байланысқан қабатқа жіберіледі [127, б.81-84, 129].

### 3.5.3 Үлгінің гиперпараметрлері

Эмпирикалық түрде әртүрлі гипер-параметрлер тексерілді. LSTM конфигурациясы үшін жасырын күй өлшемі 128-ге тең, активация функциясы ретінде ReLU қолданылған. Барлық сәулеттегі эпохалар саны 200 болып белгіленді.

Оқытудың бірнеше кезеңінен кейін біз мәтінді жіктеу моделін ала аламыз.

Ұсынылған нейрондық желі LSTM жіктеуіші арқылы бағаланады. Оңтайландыру ретінде Adam алгоритмі алынды. Жоғалту функциясы – binary\_crossentropy. Бинарлы кросс-энтропия келесі формула бойынша есептеледі:

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i * \log(p(y_i)) + (1 - y_i) * \log(1 - p(y_i)) \quad (3.15)$$

Шығыс қабатта тесттік мәліметтер экстремистік немесе бейтарап класстардың біріне жіктеледі.

### 3.5.4 Үлгіні бағалау параметрлері

Үлгі құрастырылғаннан кейін оның өнімділігін бағалау керек. Тиісті бағалау өте маңызды, өйткені жіктеуіштің дәлдігін білместен оны нақты тапсырмаларда қолдануға болмайды. Жіктеуішті бағалаудың көптеген әдістері, сонымен қатар

көптеген көрсеткіштері бар. Негізгі көрсеткіш – бұл тест жиынтығындағы дұрыс жіктелген даналардың саны, ол тест жиынтығындағы даналардың жалпы санына бөлінеді [119, б.80].

Машиналық оқыту тапсырмаларында метрикалар модельдердің сапасын бағалау және әр түрлі алгоритмдерді салыстыру үшін қолданылады. Мұнда біз жіктеу мәселелеріндегі кейбір сапа критерийлерін қарастырамыз.

*Accuracy*

Қарапайым жағдайда, мұндай метрика жіктеуіш дұрыс шешім қабылдаған құжаттардың үлесі бола алады.

$$Accuracy = \frac{P}{N} \quad (3.16)$$

мұндағы, P – жіктеуіш дұрыс шешім қабылдаған құжаттар саны, ал N - оқу таңдамасының мөлшері.

*Дәлдік және толықтық (Precision and Recall)*

Дәлдік пен толықтық – бұл ақпарат алудың көптеген алгоритмдерін бағалауда қолданылатын өлшемдер. Кейде олар өздігінен қолданылады, кейде F-шама немесе R-Precision сияқты туынды метрикаларға негіз бола алады. Дәлдік пен толықтықтың мәні өте қарапайым болып табылады.

Жүйенің дәлдігі – жүйенің берілген санатқа тағайындаған барлық құжаттарының шын мәнінде сол санатқа жататын құжаттар ішіндегі үлесі.

Жүйенің толықтығы – жіктеуіш берілген санатқа жатады деп анықтаған құжаттардың тесттік жинақтағы сол санаттың барлық құжаттарына қатынасы.

Бұл мәндерді әр санат үшін бөлек құрастырылатын контингенттік кестесі негізінде оңай есептеуге болады [129].

$$Precision = \frac{TP}{TP+FP} \quad (3.17)$$

$$Recall = \frac{TP}{TP+FN} \quad (3.18)$$

мұндағы TP – ақиқат оң нәтиже, классификатор объектіні қарастырылып отырған сыныпқа дұрыс жатқызды.

$$TP_i = T_i \quad (3.19)$$

FP – жалған оң нәтиже, классификатор объектіні қарастырылып отырған сыныпқа қате жатқызады.

$$FP_i = \sum_{c \in Classes} F_{i,c} \quad (3.20)$$

FN – жалған теріс нәтиже, классификатор объект қарастырылып отырған сыныпқа жатпайды деп қате болжам жасайды.

$$FN_i = \sum_{c \in Classes} F_{c,i} \quad (3.21)$$

TN – ақиқат теріс нәтиже, классификатор объект қарастырылып отырған сыныпқа жатпайды деп дұрыс болжам жасайды [130].

$$TN_i = All - TP_i - FP_i - FN_i \quad (3.22)$$

*Confusion Matrix (дәлсіздік матрицасы)*

Іс жүзінде дәлдік пен толықтық мәндерін дәлсіздік матрицасын пайдаланып есептеу әлдеқайда ыңғайлы (confusion matrix). Егер санаттардың саны салыстырмалы

түрде аз болса (100-150 санаттан артық емес), бұл тәсіл классификатор жұмысының нәтижелерін жеткілікті түрде айқын көрсетуге мүмкіндік береді.

Дәлсіздік матрицасы  $N$ -ға  $N$  матрица, мұндағы  $N$  – класстар саны. Бұл матрицаның бағандары сарапшылардың шешімдері үшін, ал жолдар жіктеуіштердің шешімдері үшін сақталған.

Мұндай матрица болған жағдайда әр сынып үшін дәлдік пен толықтық өте қарапайым есептеледі. Дәлдік матрицаның сәйкес диагональды элементінің және санаттың барлық жолдарының қосындысының қатынасына тең. Толықтық – матрицаның диагональды элементінің және сыныптың барлық бағанының қосындысының қатынасына тең. Жіктеуіштің алынған дәлдігі оның барлық кластары үшін дәлдігінің орташа арифметикалық мәні ретінде есептеледі. Толықтық та дәл осылай есептеледі. Техникалық тұрғыдан бұл тәсілді макроорташа деп атайды.

*F-өлшем*

Дәлдік пен толықтық неғұрлым жоғары болса, соғұрлым жақсы болатыны анық. Бірақ нақты өмірде максималды дәлдік пен толықтыққа бір уақытта қол жеткізу мүмкін емес, сондықтан белгілі бір тепе-теңдікті іздеу керек. Сондықтан, алгоритмнің дәлдігі мен толықтығы туралы ақпаратты біріктіретін қандай да бір көрсеткіш болуы керек. F-өлшем дәл осындай метрика болып табылады.

F өлшем – дәлдік пен толықтық арасындағы гармоникалық орташа мән. Егер дәлдік немесе толықтық нөлге ұмтылса, ол да нөлге ұмтылады.

$$F - Measure = 2 * \frac{Precision * Recall}{Precision + Recall} \quad (3.23)$$

Бұл формула дәлдік пен толықтыққа бірдей салмақ береді, сондықтан F өлшемі дәлдік пен толықтықтың төмендеуімен бірдей болады.

$$F - Measure = (\beta^2 + 1) * \frac{Precision * Recall}{\beta^2 * Precision + Recall} \quad (3.24)$$

Дәлдікке басымдық берілетін жағдайда мұндағы  $\beta$ ,  $0 < \beta < 1$  диапазонындағы мәндерді қабылдайды, ал  $\beta > 1$  толықтыққа басымдық береді.  $\beta = 1$  болған кезде формула алдыңғы формулаға келеді және теңдестірілген F өлшемін алынады (оны F1 деп те атайды).

*ROC қисығы астындағы аудан (AUC-ROC)*

ROC қисығы (receiver operating characteristics, қабылдағыштың жұмыс сипаттамасы) – бұл әр түрлі шекті параметрлердегі жіктеу мәселелеріне арналған өнімділікті өлшеу белгісі. ROC – бұл ықтималдық қисығы, ал AUC-ROC бөліну дәрежесін немесе өлшемін білдіреді. Бұл өлшем модельдің санаттарды қаншалықты ажырата алатындығын айтады. AUC-ROC неғұрлым жоғары болса, модель 0 классты 0 және 1 классты 1 деп болжауды неғұрлым жақсы жүргізеді [128, б.87].

### **3.6 Эксперименталды бөлім**

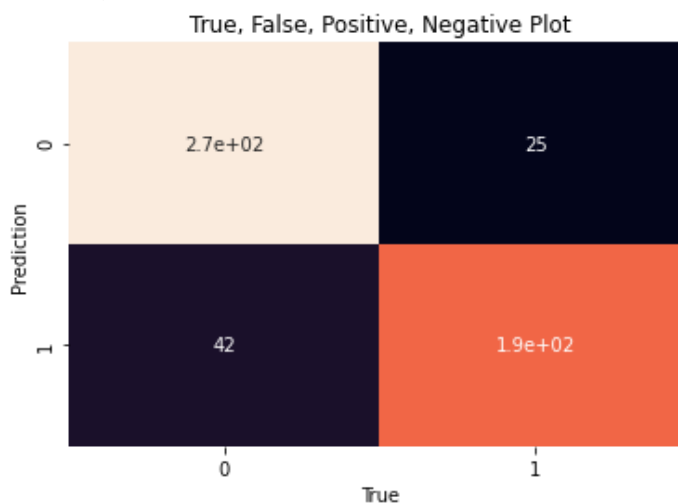
Бұл бөлімде TF-IDF+биграммдар, Word2vec+биграммдар және Bag of Words+биграммдарға негізделіп құрастырылған үш түрлі әдіс бойынша олардың ішіндегі тиімдісін анықтауға қатысты жүргізілген зерттеу нәтижелері келтіріледі [131-133].

LSTM негізіндегі терең оқыту үлгісінің кірісіне алдымен TF-IDF және биграммдар комбинациясы берілді. Эксперимент нәтижесінде мәтіндерді жіктеу дәлдігі 0.87, экстремистік мәтіндерді анықтаудың F1-шамасы 0.85 құрады (сурет 3.14).

	precision	recall	f1-score	support
0	0.86	0.91	0.89	293
1	0.88	0.82	0.85	232
accuracy			0.87	525

Сурет 3.14 – TF-IDF және биграммдар арқылы жіктеу нәтижесі

Дәлсіздік матрицасы нақты мақсатты мәндерді машиналық оқыту моделі болжайтын мәндермен салыстырады. Бұл бізге жіктеу моделінің қаншалықты жақсы жұмыс істейтіндігі және қандай қателіктер жіберетіндігі туралы көрініс береді. Мақсатты айнымалының екі мәні бар: оң немесе теріс. Бағандар мақсатты айнымалының нақты мәндерін білдіреді. Жолдар мақсатты айнымалының болжамды мәндерін білдіреді (сурет 3.15).



Сурет 3.15 – TF-IDF және биграммдар арқылы жіктеудің дәлсіздік матрицасы

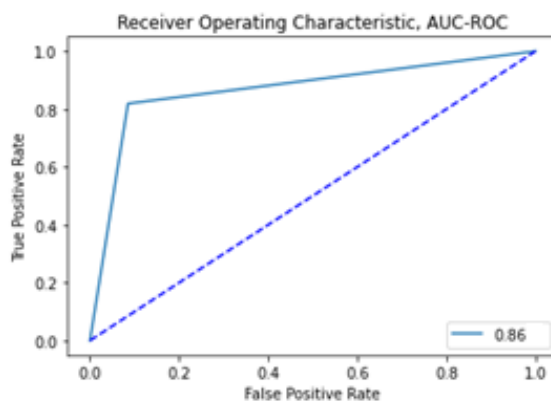
Нақты оң (TP) = 2.7e+02 (733); 560 оң класс деректері модель бойынша дұрыс жіктелген.

Нақты теріс (TN) = 1.9e+02 (516); 516 теріс класс деректері модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 25; теріс кластағы 25 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген

Жалған теріс (FN) = 42; оң кластағы 42 дерек модель бойынша теріс классқа жатады деп қате жіктелген. Осы белгілер бойынша жіктеудің AUC-ROC түріндегі нәтижесі төмендегі суретте келтірілген (сурет 3.16):





Сурет 3.16 – TF-IDF және биграммдар арқылы жіктеудің AUC-ROC мәні

AUC-ROC мәні 0.86 құрайтындығын көреміз, бұл жіктеуіштің оң сынып мәндерін теріс сынып мәндерінен ажырата алуының үлкен мүмкіндігі бар екендігін көрсетеді.

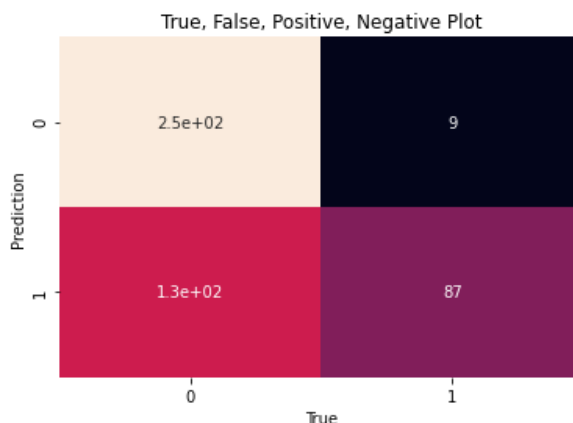
Келесі кезекте терең оқыту алгоритмдерінің кірісіне Word2vec және биграммдар комбинациясын беру арқылы келесі нәтижелер алынды (сурет 3.17):

	precision	recall	f1-score	support
0	0.66	0.96	0.78	257
1	0.91	0.40	0.56	217
accuracy			0.71	474

Сурет 3.17 – Word2vec және биграммдар арқылы жіктеу нәтижесі

Эксперимент нәтижесінде мәтіндерді жіктеу дәлдігі 0.71, экстремистік мәтіндерді анықтаудың F1-шамасы 0.56 құрады.

Word2vec және биграммдар нәтижесінде алынған дәлсіздік матрицасы төменде келтірілген (сурет 3.18).



Сурет 3.18 – Word2vec және биграммдар арқылы жіктеудің дәлсіздік матрицасы

Нақты оң (TP) = 2.5e+02 (680); оң класстың 680 дерегі модель бойынша дұрыс жіктелген.

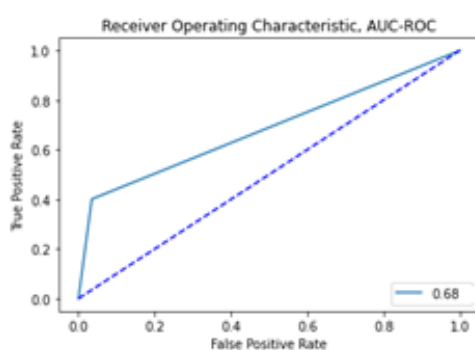
Нақты теріс (TN) = 87; 87 теріс класс деректері модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 9; теріс кластағы 9 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген

Жалған теріс (FN) =  $1.3e+02$  (353) ; оң кластағы 353 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

Матрицадан көріп тұрғанымыздай, Word2vec және биграммдарды қолдану барысында модельдегі жалған теріс нәтижелердің саны өте үлкен (353 дерекке дұрыс жіктеу жасалмаған). Берілген нәтижеге сүйеніп, аталған белгілер комбинациясын экстремистік мәтіндерді жіктеуде қолданудың тиімсіз екендігін көруге болады.

Осы белгілер бойынша жіктеудің AUC-ROC түріндегі нәтижесі төмендегі суретте келтірілген (сурет 3.19):



Сурет 3.19 – Word2vec және биграммдар арқылы жіктеудің AUC-ROC мәні

AUC-ROC шамасы 0.68-ді құрайды, бұл жіктеуіш позитивті және негативті класс нүктелерін ажыратуда жақсы нәтиже көрсетпейтіндігін білдіреді (0.5 шамасына жақын болған сайын, жіктеуіш өнімділігі төмен деп бағаланады). Жіктеуіш барлық деректер үшін кездейсоқ классты немесе тұрақты класты болжайтын болады. Жіктеуіш үшін AUC-ROC мәні неғұрлым жоғары болса, оның оң және теріс кластарды ажырата білу қабілеті соғұрлым жоғары болады. Эксперименттің келесі бөлімінде мәтіндерге терең оқыту алгоритмдері негізінде Bag of words және биграмм белгілерін қолдану арқылы жіктеу жүргізілді, аталған әдіс келесідей нәтижелер көрсетті (сурет 3.20):

	precision	recall	f1-score	support
0	0.74	0.98	0.85	293
1	0.96	0.57	0.72	232
accuracy			0.80	525

Сурет 3.20 – Bag of words және биграмм арқылы жіктеу нәтижесі

Модельдің дәлдігі 0.8, экстремистік мәтіндерді анықтаудың F1-шамасы 0.72 құрайды. Аталған белгілерді қолдану кезіндегі дәлсіздік матрицасы төмендегі суретте келтірілген (сурет 3.21):



Сурет 3.21 – Bag of words және биграмм арқылы жіктеудің дәлсіздік матрицасы

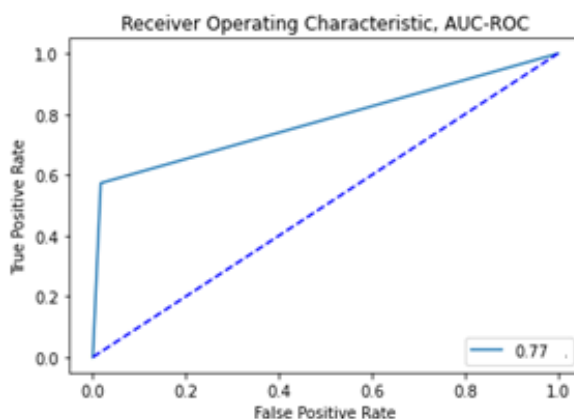
Нақты оң (TP) =  $2.9e+02$  (788); оң класстың 788 дерегі модель бойынша дұрыс жіктелген.

Нақты теріс (TN) =  $1.3e+02$  (353); теріс класстың 353 дерегі модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 5; теріс кластағы 5 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген.

Жалған теріс (FN) = 99 ; оң кластағы 99 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

Bag of words және биграммдар әдісінде да жалған теріс нәтижелердің жиі кездесетінін байқаймыз. Дегенмен, жіктеу нәтижелері бойынша белгілерді экстремистік мәтіндерді анықтау есебіне пайдалануға болады деп есептейміз. Аталған белгілер комбинациясы кезіндегі AUC-ROC нәтижесі төмендегі суретте келтірілген (сурет 3.22):



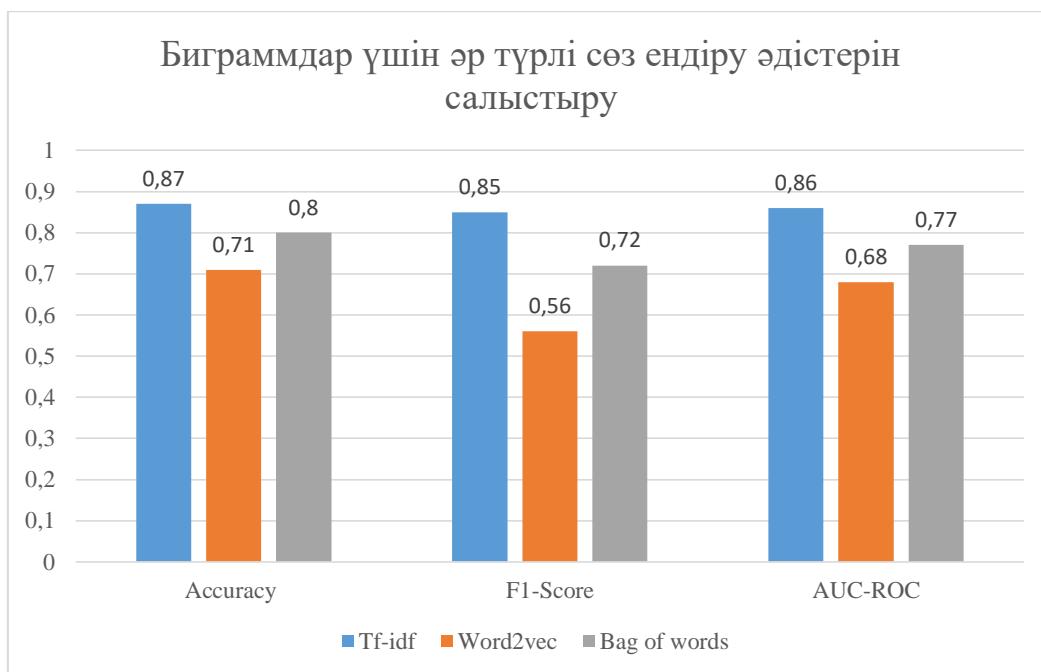
Сурет 3.22 – Bag of words және биграмм арқылы жіктеудің AUC-ROC мәні

AUC-ROC шамасы 0.77 құрайды, бұл жіктеуіштің Word2vec және биграммдар әдісімен салыстырғанда Bag of words және биграммдар әдісі негізінде жақсырақ жұмыс істейтіндігін көрсетеді.

Мәтінді TF-IDF, Word2vec, Bag of words және биграммдар әдісі негізінде жіктеудің салыстырмалы нәтижелері кестеде және диаграмма түрінде көрсетілген (кесте 3.1, сурет 3.23).

Кесте 3.1 – TF-IDF, Word2vec, Bag of words және биграммдар әдісі негізінде жіктеудің салыстырмалы нәтижелері

Терең оқыту алгоритмінде қолданылған белгілер	Accuracy	F1-Score	AUC-ROC
TF-IDF+bigram	0.87	0.85	0.86
Word2vec+bigram	0.71	0.56	0.68
Bag of words+bigram	0.80	0.72	0.77



Сурет 3.23 – TF-IDF, Word2vec, Bag of words және биграммдар әдісі негізінде жіктеудің салыстырмалы нәтижелері

Кесте 3.1-ден көріп тұрғанымыздай, TF-IDF+bigram әдісі барлық бағалау параметрлері бойынша ең жақсы нәтижені көрсетті (модель дәлдігі 0.87, экстремистік мәтіндерді анықтаудың F1-шамасы 0.85 және AUC-ROC мәні 0.86).

Эксперимент нәтижелері бойынша екінші кезектегі дәлдігі жоғары үлгі Bag of words+bigram негізіндегі үлгі болып есептеледі, модель дәлдігі 0.8, экстремистік мәтіндерді анықтаудың F1-шамасы 0.72, AUC-ROC мәні 0.77 құрайды.

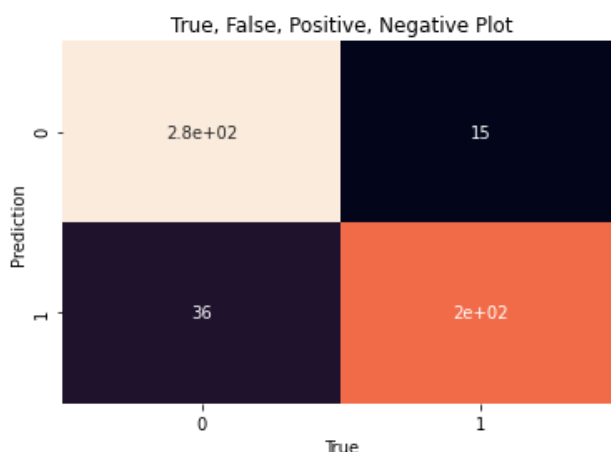
Word2vec+bigram негізіндегі үлгі берілген тапсырма үшін аса қолайлы болмайды, себебі жіктеу нәтижелері төмен көрсеткіштерді көрсетті, соның ішінде AUC-ROC мәні 0.68 құрайды, ал бұл берілген үлгінің оң және теріс класс деректерін айқын ажырату қабілетінің төменірек екендігін көрсетеді.

Аталған алгоритмдердің дәлдігін арттыру мақсатында қазақ тіліне арналған стемминг алгоритмі қолданылды. Стемминг алгоритмін қолдану нәтижесінде TF-IDF+bigram белгілерінің комбинациясы келесідей нәтиже көрсетті (сурет 3.24).

	precision	recall	f1-score	support
0	0.89	0.95	0.92	293
1	0.93	0.84	0.88	232
accuracy			0.90	525

Сурет 3.24 – Алдын ала стемминг алгоритмі орындалған биграммдарға TF-IDF әдісін қолдану арқылы алынған жіктеу нәтижесі

Яғни, стеммерсіз жіктеумен салыстырғанда үлгінің дәлдігі 0.87-ден 0.9-ға, экстремистік мәтіндерді анықтаудың F1-өлшемі 0.85-тен 0.88-ге жоғарылағандығын байқаймыз. TF-IDF+bigram+стемминг нәтижесінде алынған дәлсіздік матрицасы төменде келтірілген (сурет 3.25):



Сурет 3.25 – Алдын ала стемминг алгоритмі орындалған биграммдарға TF-IDF әдісін қолдану арқылы жіктеудің дәлсіздік матрицасы

Нақты оң (TP) =  $2.8e+02$  (761); оң класстың 761 дерегі модель бойынша дұрыс жіктелген.

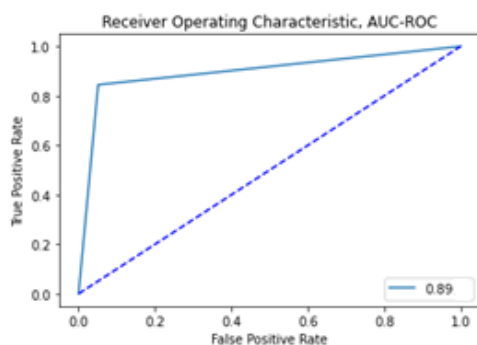
Нақты теріс (TN) =  $2e+02$  (544); теріс класстың 544 дерегі модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 15; теріс кластағы 15 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген.

Жалған теріс (FN) = 36 ; оң кластағы 36 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

Нақты оң, нақты теріс, жалған оң және жалған теріс нәтижелердің стеммерсіз жағдаймен салыстырғанда жақсарғандығын аңғарамыз.

AUC-ROC шамасы 0.89 құрайды, берілген шама стеммерсіз жағдаймен салыстырғанда 0.86-дан 0.89 шамасына дейін артқан (сурет 3.26).



Сурет 3.26 – Алдын ала стемминг алгоритмі орындалған биграмдарға TF-IDF әдісін қолдану арқылы жіктеудің AUC-ROC мәні

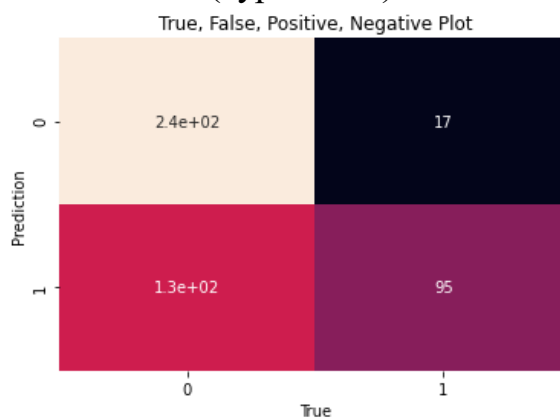
Бұл нәтиже жіктеуіштің оң және теріс класстарды ажырату өнімділігінің жоғарылағандығын білдіреді.

Келесі кезекте Word2Vec және биграмм белгілеріне стеммер қосылып, келесідей нәтиже алынды (сурет 3.27).

	precision	recall	f1-score	support
0	0.65	0.94	0.77	262
1	0.85	0.42	0.56	228
accuracy			0.69	490

Сурет 3.27 – Алдын ала стемминг алгоритмі орындалған биграмдарға Word2Vec әдісін қолдану арқылы жіктеу нәтижесі

Модельдің дәлдігі 0.69, экстремистік мәтіндерді анықтаудың F1-шамасы 0.56 құрайды. Бұл әдісте стеммерсіз жағдаймен салыстырғанда керісінше модель нәтижелерінің төмендегенін байқаймыз (сурет 3.28).



Сурет 3.28 – Алдын ала стемминг алгоритмі орындалған биграмдарға Word2Vec әдісін қолдану арқылы жіктеудің дәлсіздік матрицасы

Нақты оң (TP) = 2.4e+02 (652); оң класстың 652 дерегі модель бойынша дұрыс жіктелген.

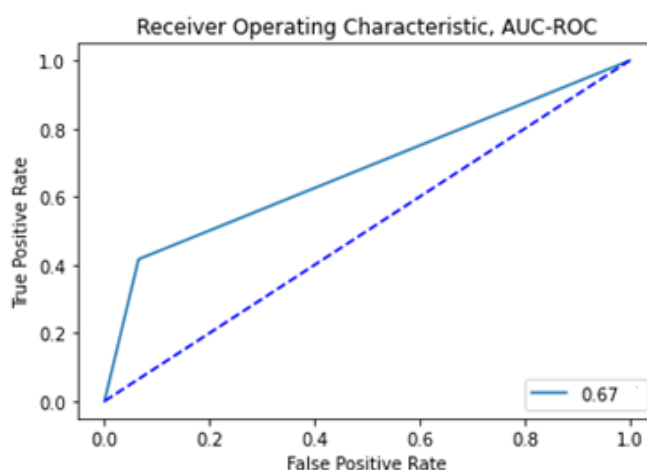
Нақты теріс (TN) = 95; теріс класстың 95 дерегі модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 17; теріс кластағы 17 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген.

Жалған теріс (FN) =  $1.3e+02$  (353) ; оң кластағы 353 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

Стеммерсіз жағдаймен салыстырғанда нақты оң, нақты теріс нәтижелердің төмендеп, жалған оң және жалған теріс нәтижелердің артқанын аңғарамыз.

AUC-ROC шамасы 0.67 құрайды, берілген нәтиже стеммер қолданылған жағдайда да Word2vec+bigram комбинациясының мәтінді жіктеуде әлі де төмен нәтиже көрсететінін білдіреді (сурет 3.29).

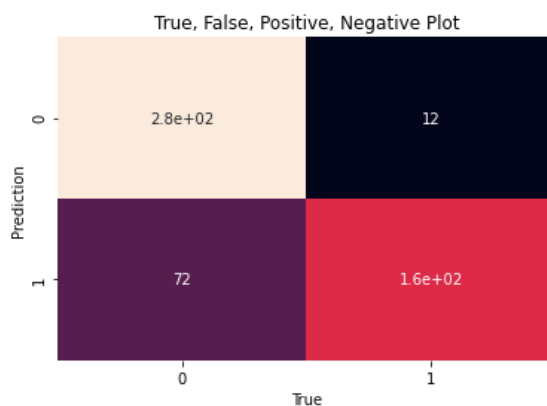


Сурет 3.29 – Алдын ала стемминг алгоритмі орындалған биграммдарға Word2vec әдісін қолдану арқылы жіктеудің AUC-ROC мәні

Келесі кезекте Bag of words+bigram белгілерімен қатар стеммерді қолдану келесідей нәтиже көрсетті. Үлгінің дәлдігі 0.80-нен 0.84-ке, экстремистік мәтіндерді анықтаудың F1-шамасы 0.72-ден 0.79-ға дейін көтерілді (кесте 3.2, сурет 3.30).

Кесте 3.2 – Алдын ала стемминг алгоритмі орындалған биграммдарға Bag of Words әдісін қолдану арқылы жіктеу нәтижесі

	precision	recall	F1-score	support
0	0.80	0.96	0.87	293
1	0.93	0.69	0.79	232
accuracy	0.84			525



Сурет 3.30 – Алдын ала стемминг алгоритмі орындалған биграммдарға Bag of Words әдісін қолдану арқылы жіктеудің дәлсіздік матрицасы

Нақты оң (TP) =  $2.8e+02$  (761); оң класстың 761 дерегі модель бойынша дұрыс жіктелген.

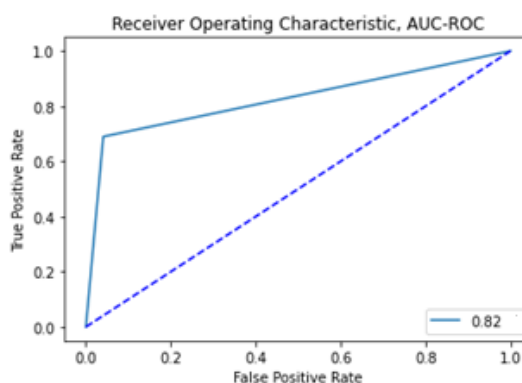
Нақты теріс (TN) =  $1.6e+02$  (435); теріс класстың 435 дерегі модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 12; теріс кластағы 12 дерек модель бойынша оң класқа жатады деп дұрыс жіктелмеген.

Жалған теріс (FN) = 72; оң кластағы 72 дерек модель бойынша теріс класқа жатады деп қате жіктелген.

Стеммерсіз Bag of words+bigram жағдайымен салыстырғанда нақты оң және жалған теріс нәтижелердің азайып, нақты теріс және жалған оң нәтижелердің артқанын байқаймыз.

AUC-ROC шамасы стеммерсіз жағдаймен салыстырғанда 0.77-ден 0.82-ге көтерілді (сурет 3.31).



Сурет 3.31 - Алдын ала стемминг алгоритмі орындалған биграммдарға Bag of Words әдісін қолдану арқылы жіктеудің AUC-ROC мәні

Берілген нәтиже стеммерді қолданған жағдайда Bag of words+bigram белгілері әдісі модельдің жұмыс өнімділігінің жақсара түсетіндігін көрсетеді.

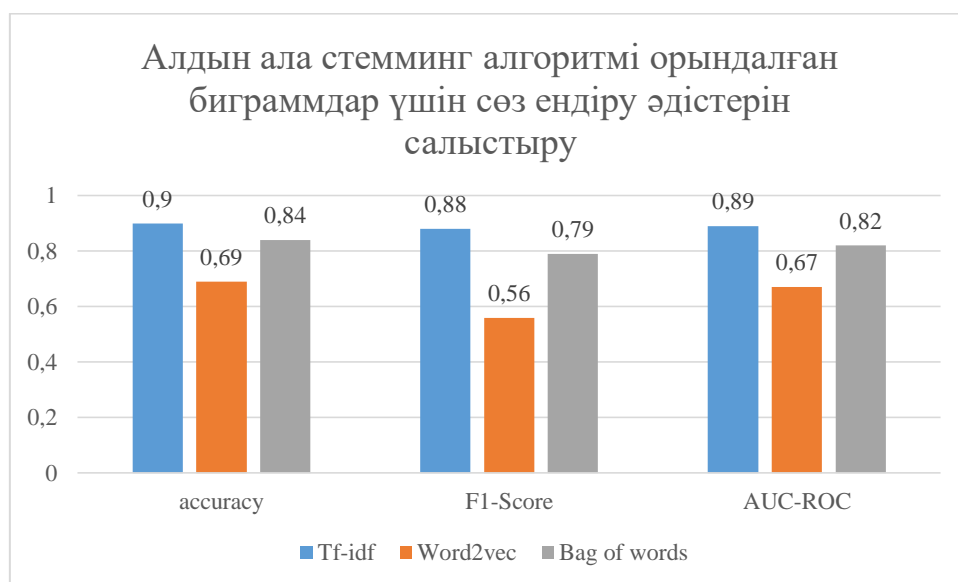
Мәтінді TF-IDF, Word2vec, Bag of words және биграммдар әдісі және стемминг алгоритмі негізінде жіктеудің салыстырмалы нәтижелері кесте 3.3-те көрсетілген.



Кесте 3.3 – Мәтінді TF-IDF, Word2vec, Bag of words және биграммдар әдісі және стемминг алгоритмі негізінде жіктеудің салыстырмалы нәтижелері

Терең оқыту алгоритмінде қолданылған белгілер	Accuracy	F1-Score	AUC-ROC
TF-IDF+bigram+стемминг	0.90	0.88	0.89
Word2vec+bigram+стемминг	0.69	0.56	0.67
Bag of words+ bigram+стемминг	0.84	0.79	0.82

Кесте 3.3-тен көріп тұрғанымыздай, TF-IDF, Word2vec, Bag of Words және биграммалармен қатар стемминг алгоритмін қолдану мәтінді терең оқыту үлгісі негізінде жіктеу есебінің дәлдігін жоғарылататындығы анықталды. Атап айтатын болсақ, стемминг алгоритмін қосу модель дәлдігін TF-IDF+bigram жағдайында 0.87-ден 0.90-ға дейін, экстремистік мәтіндерді анықтаудың F1-шама мәнін 0.85-тен 0.88-ге дейін, AUC-ROC мәнін 0.86-дан 0.89 мәніне дейін және Bag of words+ bigram жағдайында модель дәлдігін 0.80-нен 0.84-ке дейін, экстремистік мәтіндерді анықтаудың F1-шама мәнін 0.72-ден 0.79-ға дейін, AUC-ROC мәнін 0.77-ден 0.82-ге дейін арттыруға мүмкіндік берді. Ал Word2vec+bigram жағдайында керісінше модель дәлдігі 0.71-ден 0.69-ге төмендесе, экстремистік мәтіндерді анықтаудың F1-шама мәні 0.56 шамасында өзгеріссіз қалды, AUC-ROC мәні 0.68 мәнінен 0.67 мәніне дейін төмендеді (сурет 3.32).



Сурет 3.32 – Алдын ала стемминг алгоритмі орындалған биграммдар үшін сөз ендіру әдістерін салыстыру

Қорытындылай келе айтатын болсақ, экстремистік мәтіндерді анықтауға арналған үш түрлі үлгі құрылды. Ұсынылған үлгілер бойынша бірқатар эксперименттер орындалды. Жүргізілген эксперименттер нәтижесінде веб-ресурстардағы экстремистік мазмұндағы мәтіндерді анықтау үшін TF-IDF+биграммдар белгілерімен қатар стеммер алгоритмі көмегімен терең оқыту

алгоритмдеріне негізделетін үлгі ұсынылады. Аталған белгілер комбинациясы эксперименттер барысында барлық бағалау параметрлері бойынша жоғары нәтиже көрсетті.

#### 4 МАШИНАЛЫҚ ОҚЫТУ АЛГОРИТМДЕРІНЕ САЛЫСТЫРМАЛЫ ТАЛДАУ

Машиналық оқыту (ағылшынша machine learning, ML) – жасанды интеллект әдістерінің класы, оған тән белгі мәселені тікелей шешу емес, көптеген ұқсас мәселелердің шешімдерін қолдану арқылы оқыту болып табылады. Мұндай әдістерді құру үшін математикалық статистика, сандық әдістер, математикалық талдау, оңтайландыру әдістері, ықтималдықтар теориясы, графикалық теориялар, сандық түрдегі деректермен жұмыс істеудің әртүрлі әдістері қолданылады [134].

Мәліметтерді талдауға негізделген машиналық оқыту технологиясы 1950 жылы, дойбы ойынына арналған алғашқы бағдарламалар жасала бастаған кезден басталады. Компьютерлердің есептеу күшінің өсуіне байланысты олардың құратын заңдылықтары мен болжамдары бірнеше есе күрделене түсті және машиналық оқытудың көмегімен шешілетін мәселелер мен есептер ауқымы кеңейді.

Оқытудың келесідей түрі бар:

1) Прецеденттер бойынша оқыту немесе индуктивті оқыту деректердегі эмпирикалық заңдылықтарды анықтауға негізделген.

2) Дедуктивті оқыту сарапшылардың білімін формализациялауды және оларды білім базасы түрінде компьютерге беруді қамтиды.

Бақыланатын оқыту – бұл ең көп таралған жағдай. Әр жағдай «объект, жауап» жұбы болып табылады. Жауаптардың объектілердің сипаттамаларына функционалды тәуелділігін табу және объектінің сипаттамасын кіріс ретінде қабылдап, шығыста жауап қайтаратын алгоритм құру қажет болады. Сапа функционалы әдетте барлық үлгі объектілері үшін алгоритм құрған жауаптардың орташа қателігі ретінде анықталады.

– Жіктеу есебі мүмкін болатын жауаптар жиынтығының ақырлы болуымен ерекшеленеді. Оларды класс белгілері деп атайды. Класс – бұл белгі мәліметі бар барлық объектілер жиынтығы.

– Регрессия есебі жауаптың нақты сан немесе сандық вектор болуымен ерекшеленеді.

– Дәреже тағайындау (learning to rank) есебі жауаптардың бірден объектілер жиынтығында алынып, содан кейін жауаптардың мәндері бойынша сұрыпталуымен ерекшеленеді. Жіктеу немесе регрессия есептеріне дейін қысқартылуы мүмкін. Ол көбінесе ақпаратты іздеуде және мәтінді талдауда қолданылады.

– Болжау есебі (forecasting) объектілердің болашаққа болжам жасау қажет болған уақыт қатарының сегменттері болуымен ерекшеленеді. Болжау есебін шешу үшін регрессияны немесе жіктеу әдістерін бейімдеуге болады.

Бақыланбайтын оқыту. Бұл жағдайда жауаптар берілмейді және объектілер арасындағы тәуелділікті іздеу керек болады.

– Кластерлеу есебі объектілердің жұптық ұқсастығы туралы деректерді қолдану арқылы оларды кластерге топтастыруды қарастырады. Сапа функционалын әр түрлі әдіспен, мысалы, орташа кластераралық және кластерішілік арақашықтықтардың қатынасы ретінде анықтауға болады.

– Қауымдастық ережелерін іздеу есебі (association rules learning). Бастапқы мәліметтер сипаттамалар түрінде ұсынылады.

Бекіту арқылы оқыту. Нысан ретінде «жағдай, қабылданған шешім» жұптары қарастырылады, жауаптар - қабылданған шешімдердің дұрыстығын сипаттайтын функционалдық сапа мәндері (қоршаған ортаның реакциясы). болжау есебіндегідей, мұнда да уақыт факторы маңызды рөл атқарады. Қолданбалы есептердің мысалдары: инвестициялық стратегияларды қалыптастыру, технологиялық процестерді автоматты басқару, роботтарды өздігінен оқыту және т.б. [135]

Оқыту барысында пайдаланылатын мәліметтер жинағы  $A = \{A_1, A_2, \dots, A_{|A|}\}$  атрибуттар жиынтығы арқылы сипатталады, мұндағы  $|A|$  атрибуттар санын немесе  $A$  жинағының өлшемін білдіреді. Сонымен қатар, мәліметтер жинағында класс атрибуты деп аталатын арнайы мақсатты  $C$  атрибуты болады.  $C$  класс атрибуты дискретті мәндер жинағына ие, яғни  $\{c_1, c_2, \dots, c_{|C|}\}$ , мұндағы  $|C|$  - класстар саны және  $|C| \geq 2$ . Класс мәні класс белгісі деп те аталады.

Оқытуға арналған мәліметтер жинағы – реляциялық кесте. Әрбір мәліметтер жазбасы машиналық оқыту тілінде мысал, экземпляр, оқиға немесе вектор деп аталады.

Мәліметтер жинағы негізінен мысалдардан тұрады. Оқыту мақсаты –  $D$  мәліметтер жинағын ескере отырып,  $A$ -дағы атрибуттар мен  $C$ -дағы класстардың мәндерін сәйкестендіруге арналған жіктеу/болжау функциясын құру болып табылады. Аталған функция болашақта кездесетін жаңа мәліметтердің класстарының мәндерін болжау үшін пайдаланылады. Бұл функция жіктеу моделі немесе классификатор деп аталуы да мүмкін.

Оқыту үшін қолданылатын мәліметтер жиынтығы оқыту деректері (немесе оқыту жиынтығы) деп аталады. Модель оқыту алгоритмі арқылы оқыту деректері негізінде зерттелгеннен немесе құрастырылғаннан кейін, модельдің дәлдігін бағалау үшін тестілік мәліметтер жиынтығымен (немесе көрінбейтін мәліметтермен) бағаланады. Тест деректері жіктеу моделін зерттеуде пайдаланылмайды.

Тестілік деректердегі мысалдарда әдетте класс белгілері болады. Оқыту және тестілеу үшін қол жетімді мәліметтер (класстары бар) әдетте екі қиылыспайтын ішкі жиындарға бөлінеді: оқыту жиынтығы (оқыту үшін) және тест жиынтығы (тестілеу үшін) [119, б.79].

Диссертациялық жұмыстың мақсатына сай дәстүрлі машиналық оқыту әдістерін қолдану арқылы эксперименттер жүргізілді. Веб-ресурстардағы экстремистік мәтіндерді анықтау үшін шешім ағашы, мультиномиялық аңқау Байес, кездейсоқ орман, сызықтық регрессия және тірек векторлар машинасы сияқты дәстүрлі әдістер қолданылып, олардың жіктеу нәтижелері салыстырылды.

#### **4.1 Тірек векторлар машинасы**

Тірек векторлар машинасы – жіктеу және регрессиялық талдау мәселелерін шешуге қажетті алгоритмдердің жиынтығы.  $N$  өлшемді кеңістіктегі нысан екі класстың біріне жататындығына сүйене отырып, тірек векторлық машинасы барлық объектілер екі топтың біреуінде болатындай етіп ( $N - 1$ ) өлшемі бар гиперкеңістік

тұрғызады. Тірек векторлар машинасын қолдану барысында келесі екі шектеу қанағаттандырылуы керек:

$$wx_i - b \geq +1, \text{ егер } y_i = +1 \quad (4.1)$$

$$wx_i - b \leq -1, \text{ егер } y_i = -1 \quad (4.2)$$

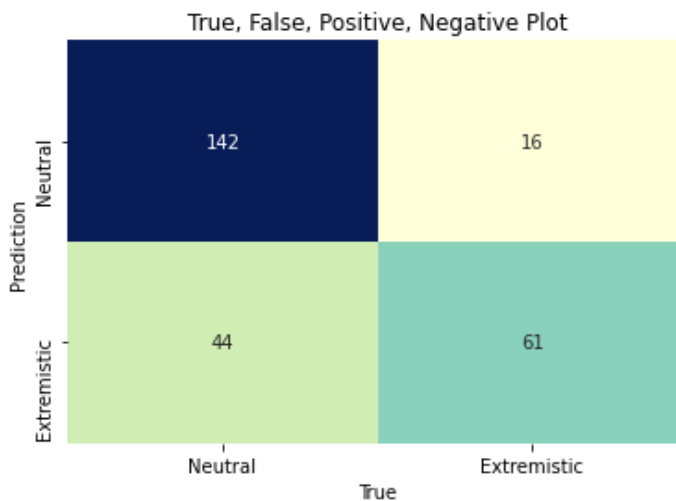
Гиперкеңістік әр класстың ең жақын мәліметтерінен бірдей қашықтыққа алыстауы үшін  $\|w\|$  мәні минимизациялануы керек, оның минимизациясы  $\frac{1}{2} \|w\|^2$  минимизациясына пара-пар. Тірек векторлар машинасындағы оңтайландыру есебі келесі сипатта болады [129, б.53]:

$$\min \frac{1}{2} \|w\|^2, y_i(x_i w - b) - 1 \geq 0, i = 1, \dots, N \quad (4.3)$$

Тәжірибе үшін берілген алгоритмді қолдану барысында келесідей нәтижелер алынды (кесте 4.1, сурет 4.1):

Кесте 4.1 – Тірек векторлар машинасы әдісі арқылы жіктеу нәтижесі

Accuracy	0,73
Precision	0,99
Recall	0,36
F1	0,53
AUC-ROC	0,68



Сурет 4.1 – Тірек векторлар машинасы арқылы жіктеу нәтижесі

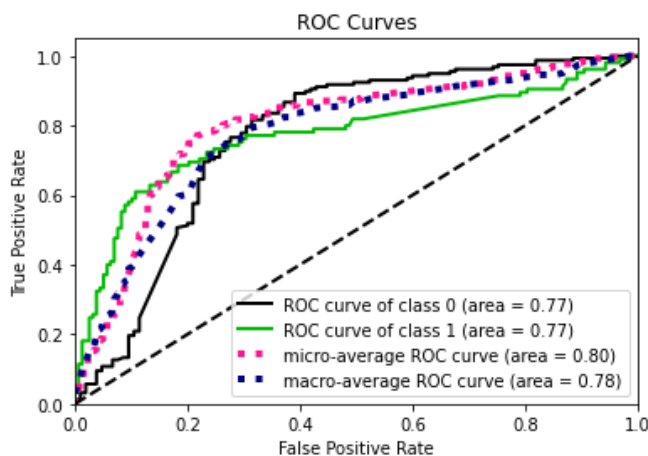
Нақты оң (TP) = 142; 142 оң класс деректері модель бойынша дұрыс жіктелген.

Нақты теріс (TN) = 61; 61 теріс класс деректері модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 16; теріс кластағы 16 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген.

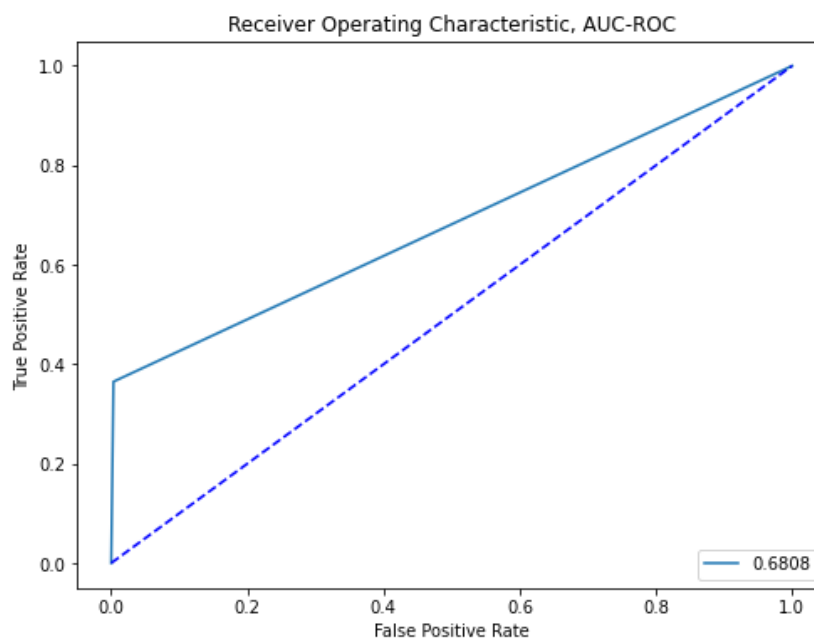
Жалған теріс (FN) = 44; оң кластағы 44 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

Осы белгілер бойынша жіктеудің ROC қисық түріндегі нәтижесі төмендегі суретте келтірілген (сурет 4.2):



Сурет 4.2 – Тірек векторлар машинасы арқылы жіктеудің ROC қисық түріндегі нәтижесі

ROC мәні 0.77 құрайтындығын көреміз, бұл жіктеуіштің оң сынып мәндерін теріс сынып мәндерінен ажырата алу мүмкіндігінің жақсы екендігін көрсетеді. AUC-ROC мәні келесі сурет 4.3-те көрсетілген:



Сурет 4.3 – Тірек векторлар машинасы арқылы жіктеудің AUC-ROC мәні

#### 4.2 Шешім ағашы

Шешім ағашы (жіктеу ағашы немесе регрессия ағашы деп те аталады) – шешім қабылдау үшін пайдаланылатын ациклді граф. Графтың әрбір бұтақты түйінінде белгілер векторындағы  $j$ -шы белгі зерттеледі. Егер белгінің мәні берілген шектен төмен болса, сол жақ бұтақ, кері жағдайда оң жақ бұтақ таңдалады. Жапырақ түйініне жеткен кезде берілген дана тиісті болатын класс жайлы шешім қабылданады.

Мәліметтерді талдау барысында шешімдер ағашын математикалық және есептеу әдістері ретінде мәліметтер жинағын сипаттау, жіктеу және жалпыландыру үшін пайдалануға және келесі түрде жазуға болады [129, б.49]:

$$(x, Y) = (x_1, x_2, x_3 \dots, x_k, Y) \quad (4.4)$$

Тәуелді Y айнымалысы талдануы және жіктелуі тиіс мақсатты айнымалы болып табылады. x векторы  $x_1, x_2, x_3$  және т.с.с. кіріс айнымалыларынан тұрады.

Шешімдер ағашы арқылы талдау кезінде бәсекелес баламалардың күтілетін мәндерін (немесе күтілетін пайдасын) есептеу үшін шешімдерді қолдаудың визуалды және аналитикалық құралы қолданылады.

Берілген алгоритмдегі оңтайландыру критерийі орташа логарифмдік шындыққа ұқсастық болып табылады:

$$\frac{1}{N} \sum_{i=1}^N [y_i \ln f_{ID3}(x_i) + (1 - y_i) \ln(1 - f_{ID3}(x_i))] \quad (4.5)$$

мұндағы  $f_{ID3}$  - шешімдер ағашы.

S белгіленген мәліметтер жиынтығы болсын. Бастапқыда шешімдер ағашында барлық мәліметтерді қамтитын жалғыз түйін болады:  $S \stackrel{\text{def}}{=} \{(x_i, y_i)\}$

$f_{ID3}^S$  тұрақты моделі келесідей анықталады:

$$f_{ID3}^S = \frac{1}{|S|} \sum_{(x,y) \in S} y \quad (4.6)$$

$f_{ID3}^S(x)$  моделі кез келген x үшін бірдей болжам қайтарады. Әрі қарай  $j=1, \dots, D$  белгілері мен t шектері ізделеді және S жиыны екі ішкі жиынға бөлінеді:

$$S_- \stackrel{\text{def}}{=} \{(x, y) | x, y \in S, x^{(j)} < t\} \quad (4.7)$$

$$S_+ \stackrel{\text{def}}{=} \{(x, y) | x, y \in S, x^{(j)} \geq t\} \quad (4.8)$$

Берілген екі ішкі жиын жаңа жапырақтық түйіндерді құрайды. Модель нәтижесінің дұрыстығы энтропия арқылы бағаланады. Энтропия – кездейсоқ шаманың белгісіздік шамасы.

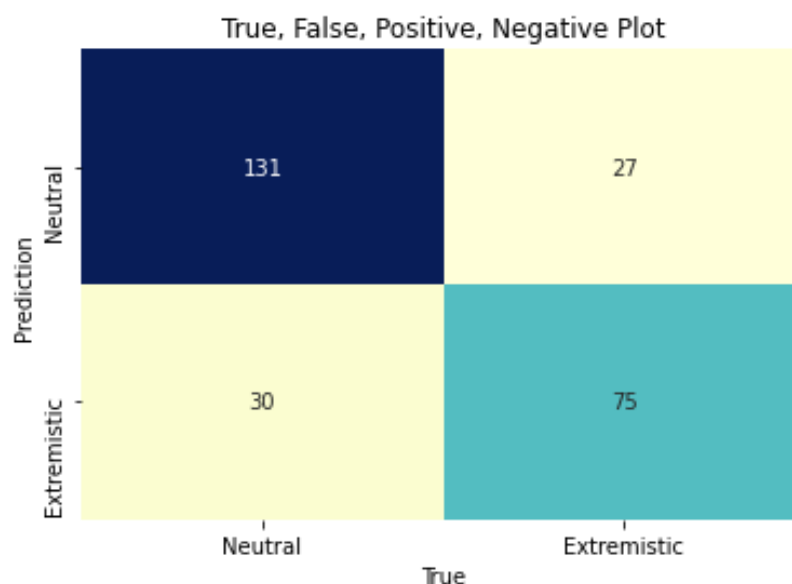
Барлық кездейсоқ шама мәндерінің ықтималдықтары тең болған кезде энтропия өз максимумына жетеді және кездейсоқ шама тек бір мәнге ғана ие болған кезде энтропия өз минимумына жетеді. S жиынының энтропиясы келесідей анықталады [129, б.50]:

$$H(S) \stackrel{\text{def}}{=} -f_{ID3}^S \ln f_{ID3}^S - (1 - f_{ID3}^S) \ln(1 - f_{ID3}^S) \quad (4.9)$$

Мәтіндерді экстремистік және бейтарап класстарға бөлу есебі үшін шешімдер ағашы әдісін қолдану барысында келесідей нәтижелер алынды (кесте 4.2, сурет 4.4):

Кесте 4.2 – Шешім ағашы әдісі арқылы жіктеу нәтижесі

Accuracy	0,77
Precision	0,95
Recall	0,49
F1	0,64
AUC-ROC	0,73



Сурет 4.4 – Шешім ағашы әдісі арқылы жіктеу нәтижесі

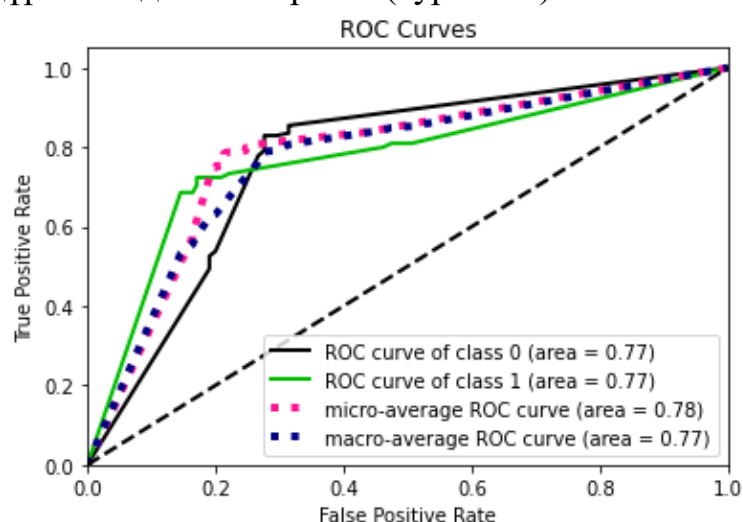
Нақты оң (TP) = 131; 131 оң класс деректері модель бойынша дұрыс жіктелген.

Нақты теріс (TN) = 75; 75 теріс класс деректері модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 27; теріс кластағы 27 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген.

Жалған теріс (FN) = 30; оң кластағы 30 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

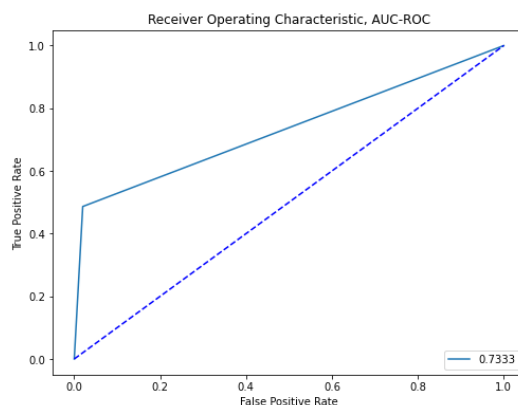
ROC мәні 0.77 құрайтындығын көреміз (сурет 4.5):



Сурет 4.5 – Шешім ағашы әдісі арқылы жіктеудің ROC қисық түріндегі нәтижесі

AUC-ROC шамасы төмендегі сурет 4.6-да көрсетілген.





Сурет 4.6 – Шешім ағашы әдісі арқылы жіктеудің AUC-ROC мәні

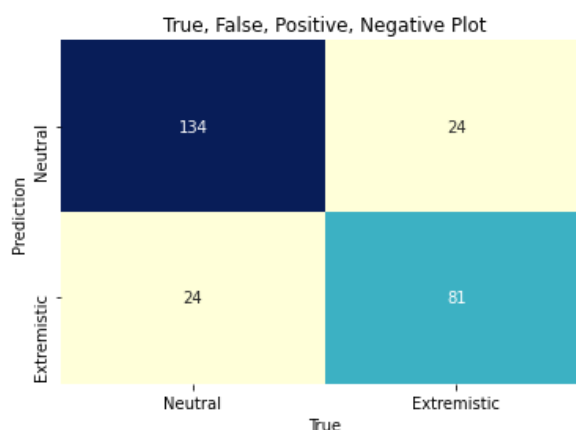
### 4.3 Кездейсоқ орман

Кездейсоқ орман – бұл Лео Брейман және Адель Катлер ұсынған машиналық оқыту алгоритмі, берілген алгоритм шешімдер ағашы ансамблін пайдаланады. Алгоритм екі негізгі идеяны қамтиды: Брейманның бэггинг әдісі және Тин Кам Хо ұсынған кездейсоқ ішкі кеңістіктер әдісі. Алгоритмнің негізгі идеясы әрқайсысы өте аз жіктеу дәлдігін беретін, бірақ олардың санының көп болуы арқасында нәтиженің жақсы болуына әсер ететін шешімдер ағашының өте көп ансамблін пайдалануға негізделеді.

Кездейсоқ орман әдісін пайдалану жағдайында келесідей нәтижелер алынды (кесте 4.3, сурет 4.7):

Кесте 4.3 – Кездейсоқ орман әдісін пайдалану арқылы жіктеу нәтижелері

Accuracy	0,73
Precision	1,0
Recall	0,36
F1	0,52
AUC-ROC	0,68



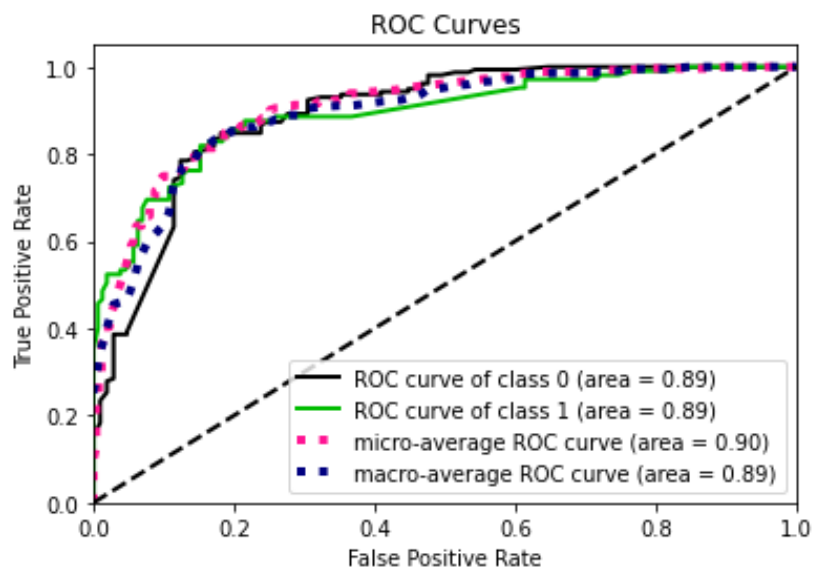
Сурет 4.7 – Кездейсоқ орман әдісі арқылы жіктеу нәтижесі  
 Нақты оң (TP) = 134; 134 оң класс деректері модель бойынша дұрыс жіктелген.

Нақты теріс (TN) = 81; 81 теріс класс деректері модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 24; теріс кластағы 24 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген.

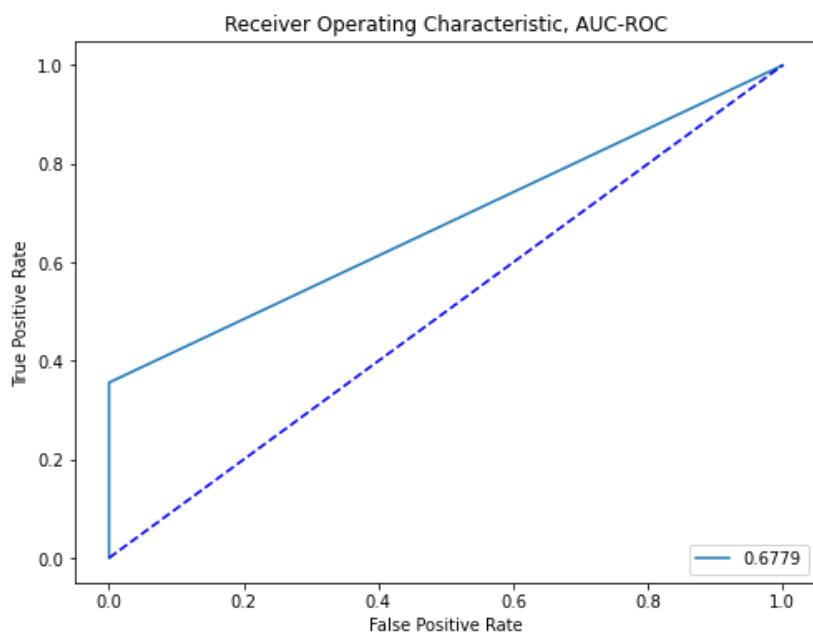
Жалған теріс (FN) = 24; оң кластағы 24 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

Осы белгілер бойынша жіктеудің ROC қисық түріндегі нәтижесі сурет 4.8-де келтірілген:



Сурет 4.8 – Кездейсоқ орман әдісі арқылы жіктеудің ROC қисық түріндегі нәтижесі

AUC-ROC шамасы төмендегі сурет 4.9-да көрсетілген.



Сурет 4.9 – Кездейсоқ орман әдісі арқылы жіктеудің AUC-ROC мәні

#### 4.4 k жақын көрші алгоритмі (k-Nearest Neighbors)

k жақын көрші алгоритмі (k-Nearest Neighbors, k-NN) – оқытудың параметрлік емес алгоритмі. Басқа машиналық оқыту алгоритмдерінен ерекшелік ретінде берілген алгоритм барлық оқыту мысалдарын жадыда сақтайды. Бұрын кездеспеген жаңа x данасы пайда болған кезде k-NN алгоритмі x-ке ең жақын k оқыту мәліметтерін табады және ең жиі кездесетін белгіні қайтарады. Екі мәліметтің жақындығы ара қашықтық функциясы арқылы анықталады. Әдетте евклидтік ара қашықтық немесе теріс косинустық ұқсастық пайдаланылады. Косинустық ұқсастық келесі формула бойынша есептеледі [129, б.57]:

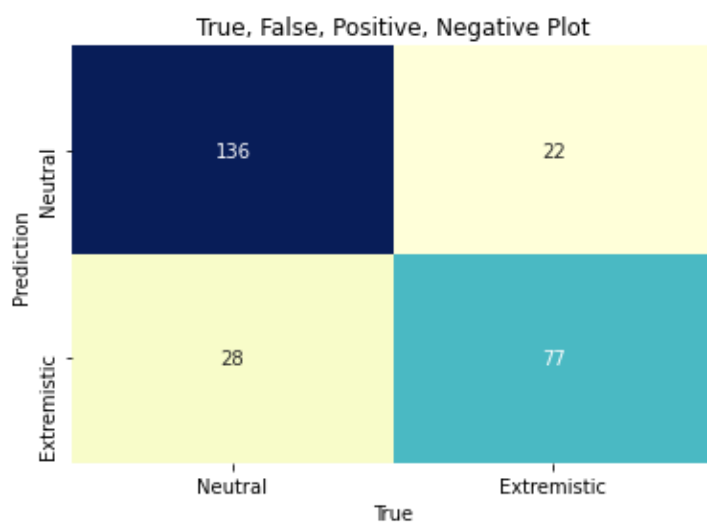
$$s(x_i, x_k) \stackrel{\text{def}}{=} \cos(\angle(x_i, x_k)) = \frac{\sum_{j=1}^D x_i^{(j)} x_k^{(j)}}{\sqrt{\sum_{j=1}^D (x_i^{(j)})^2} \sqrt{\sum_{j=1}^D (x_k^{(j)})^2}} \quad (4.10)$$

Егер векторлар арасындағы бұрыш 0 градусқа тең болса, онда бұл олардың бағыттас екендігін білдіреді және косинустық ұқсастық 1-ге тең болады. Егер векторлар ортогональ болса, косинустық ұқсастық 0-ге тең. Кері бағыттас векторлардың косинустық ұқсастығы -1-ге тең [116].

Экстремистік және бейтарап мәтіндерді жіктеу үшін берілген әдісті қолдану барысында келесідей нәтижелер алынды (кесте 4.4, сурет 4.10):

Кесте 4.4 – k жақын көрші әдісі арқылы жіктеу нәтижесі

Accuracy	0,78
Precision	0,96
Recall	0,49
F1	0,66
AUC-ROC	0,5



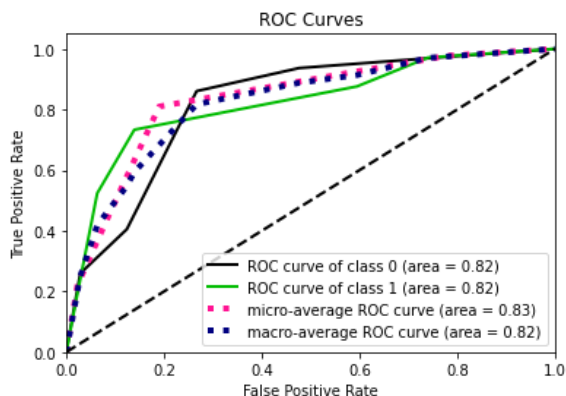
Сурет 4.10 – k жақын көрші әдісі арқылы жіктеудің дәлсіздік матрицасы

Нақты оң (TP) = 136; 136 оң класс деректері модель бойынша дұрыс жіктелген. Нақты теріс (TN) = 77; 77 теріс класс деректері модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 22; теріс кластағы 22 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген.

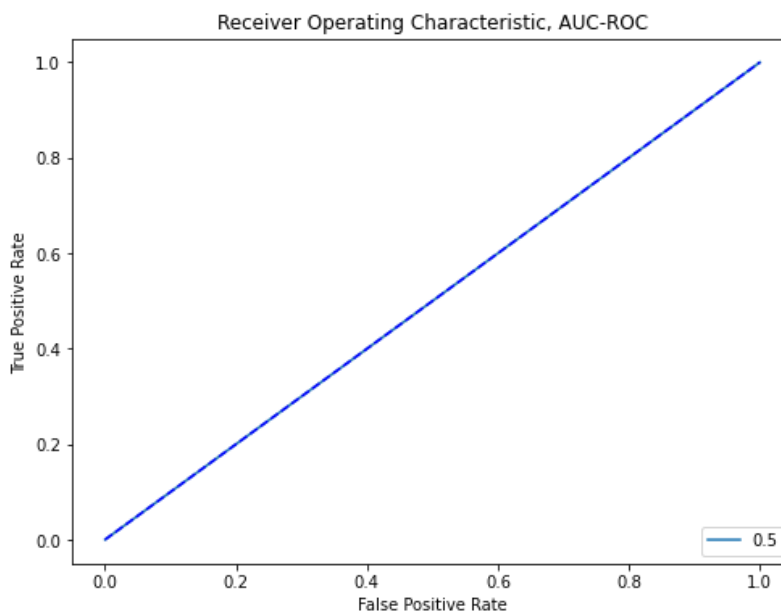
Жалған теріс (FN) = 28; оң кластағы 28 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

Осы белгілер бойынша жіктеудің ROC қисық түріндегі нәтижесі сурет 4.11-де келтірілген:



Сурет 4.11 – k жақын көрші әдісі арқылы жіктеудің ROC қисық түріндегі нәтижесі

AUC-ROC шамасы сурет 4.12-де көрсетілген.



Сурет 4.12 – k жақын көрші әдісі арқылы жіктеудің AUC-ROC мәні

#### 4.5 Аңқау Байес классификаторы

Аңқау Байес классификаторы – Байес теоремасын тәуелсіздік жайлы қатаң ұсыныстармен қатар қолдануға негізделетін қарапайым ықтималдық классификаторы. Берілген алгоритмнің артықшылығы – оқытуға, параметрлерді бағалауға және жіктеуге қажетті мәліметтер көлемінің аздығы. Классификатордың ықтималдық үлгісі – келесі шартты үлгі болып табылады,

$$p(C | F_1, \dots, F_n) \quad (4.11)$$

мұндағы  $C$  - класс,  $F_1, \dots, F_n$  - айнымалылар. Байес теоремасын қолдану арқылы келесі теңдеуді алуға болады:

$$p(C | F_1, \dots, F_n) = \frac{p(C)p(F_1, \dots, F_n | C)}{p(F_1, \dots, F_n)} \quad (4.12)$$

Тәжірибеде бөлшектің алымы қызықтырады, себебі бөлімі  $C$ -ға тәуелді емес және  $F_i$  мәндері берілген, сол себепті бөлшектің бөлімі – тұрақты шама.

Бөлшектің алымы  $p(C | F_1, \dots, F_n)$  біріккен ықтималдық үлгісіне пара-пар және оны шартты ықтималдықтың анықтамаларын пайдалану арқылы келесі түрде жазуға болады:

$$p(C | F_1, \dots, F_n) = p(C)p(F_1, \dots, F_n | C) = p(C)p(F_1 | C) p(F_2, \dots, F_n | C, F_1) = p(C)p(F_1 | C)p(F_2 | C, F_1)p(F_3, \dots, F_n | C, F_1, F_2) = p(C)p(F_1 | C)p(F_2 | C, F_1) * \dots * p(F_n | C, F_1, F_2, F_3, \dots, F_{n-1}) \quad (4.13)$$

және т.с.с.

Әрі қарай шартты ықтималдықтың «аңқау» ұсыныстарын пайдалануға болады: әрбір  $F_i$  қасиеті кез келген басқа  $F_j, i \neq j$  қасиетінен шартты тәуелсіз деп есептейік, бұл

$p(F_i | C, F_j) = p(F_i | C)$  екендігін білдіреді, осылайша біріккен үлгіні келесі түрде өрнектеуге болады:

$$p(C, F_1, \dots, F_n) = p(C)p(F_1 | C)p(F_2 | C)p(F_3 | C) * \dots * p(F_n | C) = p(C) \prod_{i=1}^n p(F_i | C) \quad (4.14)$$

Бұл тәуелсіздік жайлы ұсыныстарды,  $C$  класстық айнымалысы бойынша шартты үлестірімді келесі түрде өрнектеуге болатындығын білдіреді:

$$p(C | F_1, \dots, F_n) = \frac{1}{Z} p(C) \prod_{i=1}^n p(F_i | C) \quad (4.15)$$

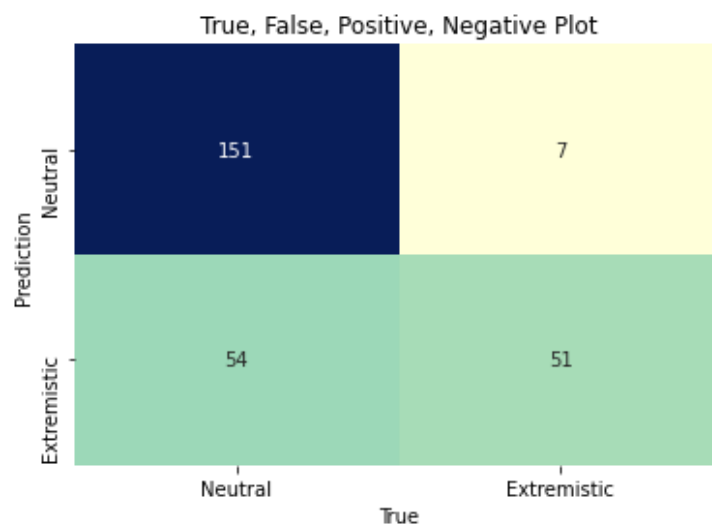
мұндағы  $Z = p(F_1, \dots, F_n)$  – бұл тек  $F_1, \dots, F_n$  –ге тәуелді масштабтық көбейткіш, яғни айнымалылардың мәндері белгілі болған жағдайда тұрақты шама болып табылады. Ықтималдық үлгісі бойынша классификатор құрылады [119, б.100, 137]:

$$classify(f_1, \dots, f_n) = \underset{c}{\operatorname{argmax}} p(C = c) \prod_{i=1}^n p(F_i = f_i | C = c) \quad (4.16)$$

Аңқау Байес классификаторы бойынша мәтіндерді экстремистік және бейтарап санатқа жіктеу барысында келесідей нәтижелер алынды (кесте 4.5, сурет 4.13):

Кесте 4.5 – Аңқау Байес әдісі арқылы жіктеу нәтижесі

Accuracy	0,86
Precision	0,80
Recall	0,89
F1	0,84
AUC-ROC	0,86



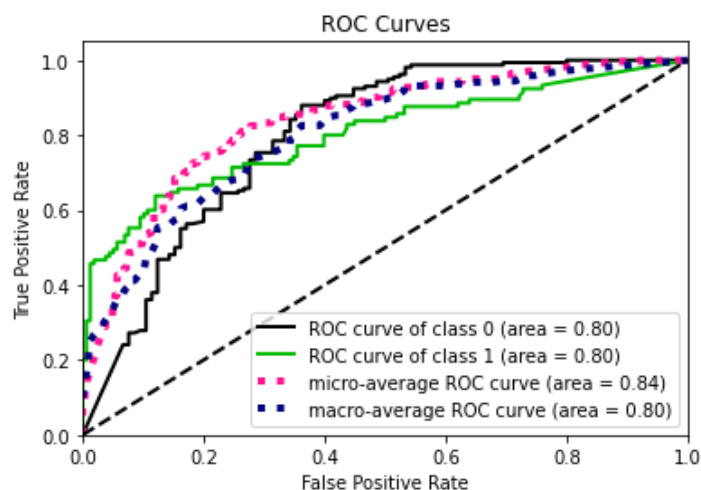
Сурет 4.13 – Аңқау Байес әдісі арқылы жіктеудің дәлсіздік матрицасы

Нақты оң (TP) = 151; 151 оң класс деректері модель бойынша дұрыс жіктелген. Нақты теріс (TN) = 51; 51 теріс класс деректері модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 7; теріс кластағы 7 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген.

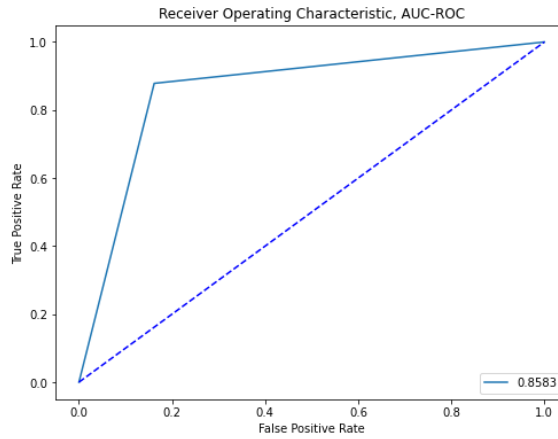
Жалған теріс (FN) = 54; оң кластағы 54 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

Аңқау Байес алгоритмі бойынша жіктеу барысында төмендегідей ROC қисық алынды (сурет 4.14):



Сурет 4.14 – Аңқау Байес әдісі арқылы жіктеудің ROC қисық түріндегі нәтижесі

AUC-ROC шамасы сурет 4.15-те көрсетілген.



Сурет 4.15 – Аңқау Байес әдісі арқылы жіктеудің AUC-ROC мәні

#### 4.6 Логистикалық регрессия

**Логистикалық регрессия** –  $y_i$ -ді  $x_i$ -ге тәуелді сызықты функция түрінде модельдеуге мүмкіндік беретін машиналық оқыту әдісі. Белгілердің  $w x_i + b$  түріндегі сызықты комбинациясы – теріс шексіздіктен оң шексіздікке дейін созылатын функция, ал  $y_i$  тек екі мәнді біріне ғана ие бола алады. Егер модельдің  $x$  данасы үшін қайтаратын мәні 0-ге жуық болса, онда оған теріс белгі тағайындалады; кері жағдайда берілген данаға оң белгі тағайындалады. Мұндай қасиеттерге ие функциялардың бірі стандартты логистикалық функция (логистикалық сигмоид) болып табылады:

$$f(x) = \frac{1}{1+e^{-x}} \quad (4.17)$$

мұндағы  $e$  – натурал логарифм негізі. Логистикалық регрессия моделі келесідей болады:

$$f_{w,b}(x) = \frac{1}{1+e^{-(wx+b)}} \quad (4.18)$$

Егер  $w$  және  $b$  мәндері сәйкесінше оңтайландырылатын болса, онда  $f(x)$  нәтижесін  $y_i$ -тің оң мәнге ие болу ықтималдығы ретінде ұсынуға болады.

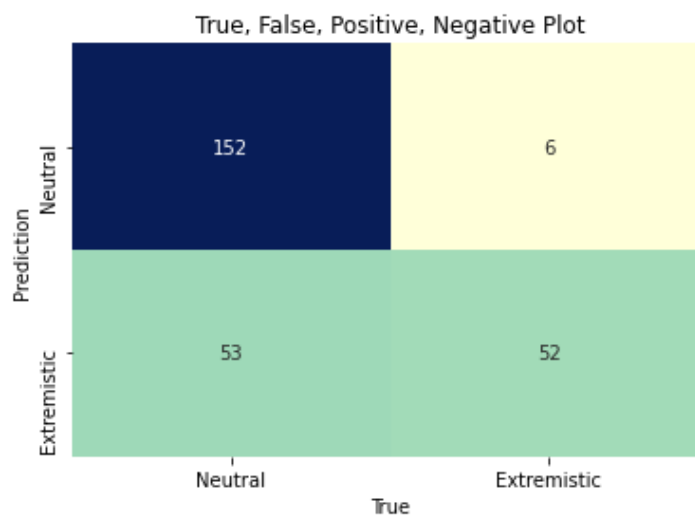
Логистикалық регрессия моделіндегі оңтайландыру критерийі максималды шындыққа ұқсастық (maximum likelihood) деп аталады. Оқытушы мәліметтердің шындыққа ұқсастығы модельге сәйкес максималдандырылады [129, б.46]:

$$L_{w,b} \stackrel{\text{def}}{=} \prod_{i=1 \dots N} f_{w,b}(x_i)^{y_i} (1 - f_{w,b}(x_i))^{(1-y_i)} \quad (4.19)$$

Веб-ресурстардағы қазақ тіліндегі экстремистік мәліметтерді анықтау тәжірибелерінде логистикалық регрессия әдісін қолдану нәтижелері кесте 4.6 мен сурет 4.16-да көрсетілген:

Кесте 4.6 – Логистикалық регрессия әдісі арқылы жіктеу нәтижесі

Accuracy	0,77
Precision	0,98
Recall	0,46
F1	0,63
AUC-ROC	0,86



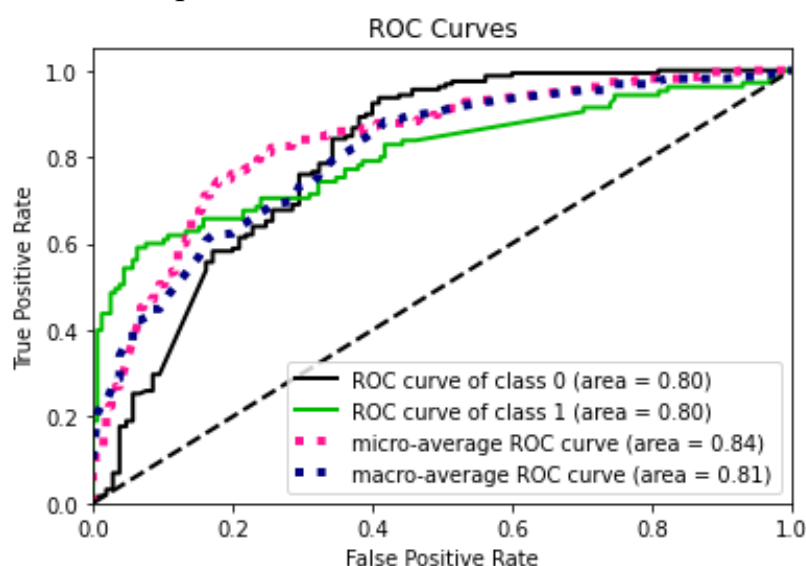
Сурет 4.17 – Логистикалық регрессия әдісі арқылы жіктеудің дәлсіздік матрицасы

Нақты оң (TP) = 152; 152 оң класс деректері модель бойынша дұрыс жіктелген. Нақты теріс (TN) = 52; 52 теріс класс деректері модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 6; теріс кластағы 6 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген.

Жалған теріс (FN) = 53; оң кластағы 53 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

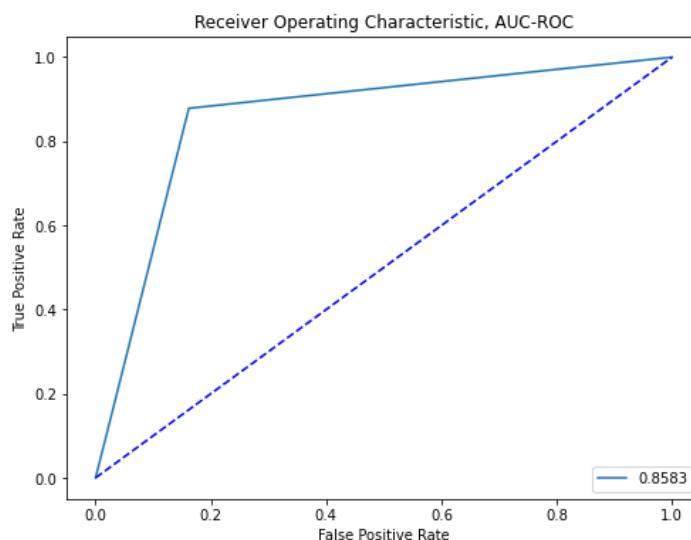
Мәтіндерді логистикалық регрессия бойынша жіктеу нәтижесінде алынған ROC қисық сурет 4.18-де келтірілген:



Сурет 4.18 – Логистикалық регрессия әдісі арқылы жіктеудің ROC қисық түріндегі нәтижесі

AUC-ROC шамасы сурет 4.19-да көрсетілген.





Сурет 4.19 – Логистикалық регрессия әдісі арқылы жіктеудің AUC-ROC мәні

#### 4.7 Градиентті бустинг

Градиентті бустинг – бұл жіктеу мен регрессияның есептерін шешуге арналған, әдетте шешімдер ағашы сияқты әлсіз болжаушы үлгілер ансамблінен тұратын болжам үлгісін құрастыратын машиналық оқыту техникасы. Кез келген оқытушымен оқытылатын алгоритмнің мақсаты – жоғалту функциясын анықтау және оны минимизациялау. Жоғалту функциясы ретінде орташа квадраттық қате таңдалған жағдайда

$$Loss = \sum (y_i - y_i^p)^2 \quad (4.20)$$

мұндағы  $y_i$ -ші мақсатты шама,  $y_i^p$ -ші болжам,  $L(y_i, y_i^p)$  – жоғалту функциясы.

Градиенттік жылжуды қолдана отырып, оқыту жылдамдығына негізделетін мәндерді жаңарту арқылы  $MSE$  минималды болатын мәндер ізделеді.

$$y_i^p = y_i^p + \alpha * \delta \sum \frac{(y_i - y_i^p)^2}{\delta y_i^p} \quad (4.21)$$

Берілген өрнек келесі түрге енеді:

$$y_i^p = y_i^p - \alpha * 2 \sum (y_i - y_i^p) \quad (4.22)$$

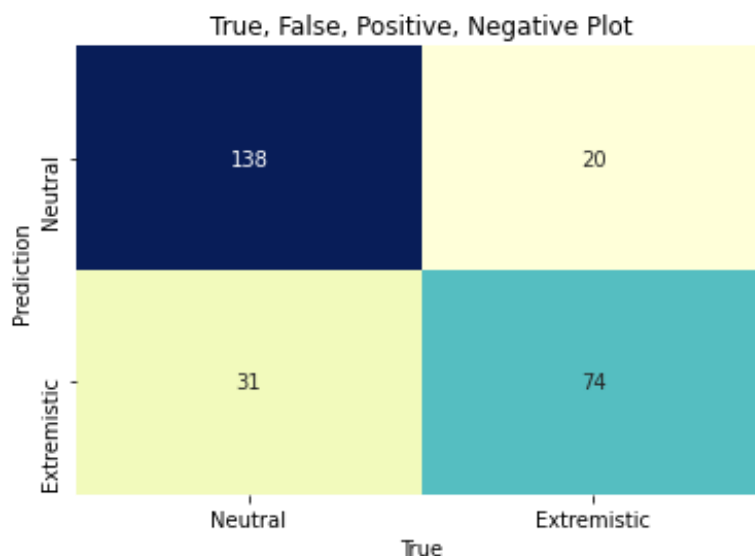
мұндағы  $\alpha$ -оқыту жылдамдығы және  $(y_i - y_i^p)$  – қалдықтар.

Осылайша, болжамдар ауытқулардың сомасы нөлге ұмтылып, болжанатын мәндер нақты мәндерге жақын болғанға дейін жаңартылып отырады [137].

Берілген алгоритмді экстремистік және бейтарап мәтіндерді жіктеу есебіне қолдану нәтижесі кесте 4.7-де, сурет 4.20-да келтірілген.

Кесте 4.7 – Градиентті бустинг әдісі арқылы жіктеу нәтижесі

Accuracy	0,75
Precision	0,99
Recall	0,41
F1	0,58
AUC-ROC	0,71



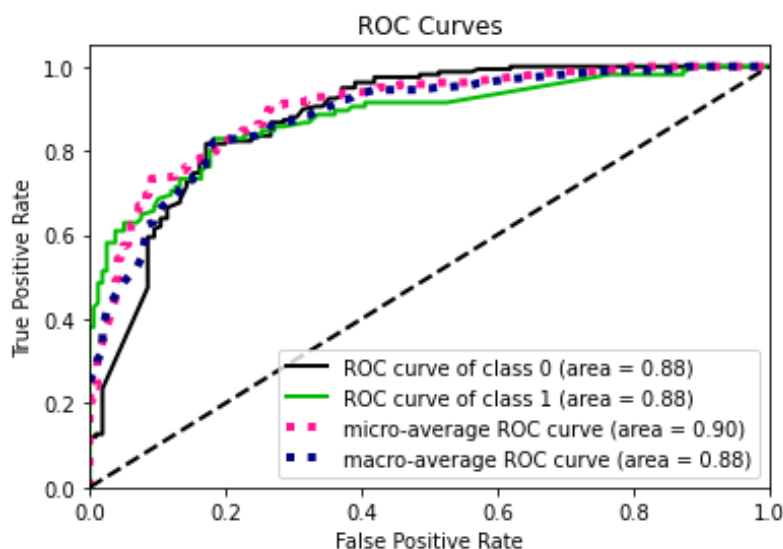
Сурет 4.20 – Градиентті бустинг әдісі арқылы жіктеудің дәлсіздік матрицасы

Нақты оң (TP) = 138; 138 оң класс деректері модель бойынша дұрыс жіктелген. Нақты теріс (TN) = 74; 74 теріс класс деректері модель бойынша дұрыс жіктелген.

Жалған оң (FP) = 20; теріс кластағы 20 дерек модель бойынша оң классқа жатады деп дұрыс жіктелмеген.

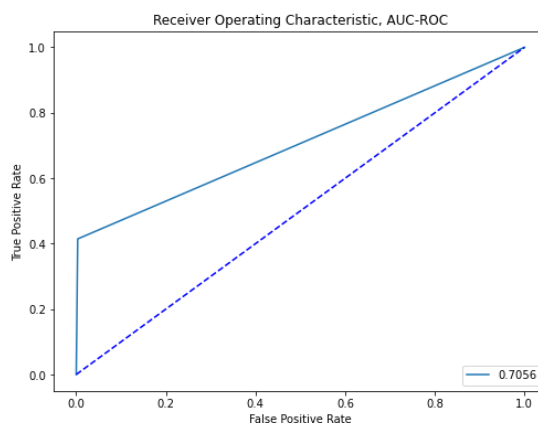
Жалған теріс (FN) = 31; оң кластағы 31 дерек модель бойынша теріс классқа жатады деп қате жіктелген.

Градиентті бустинг алгоритмі бойынша жіктеу барысында сурет 4.21-дегідей ROC қисық алынды:



Сурет 4.21 – Градиентті бустинг әдісі арқылы жіктеудің ROC қисық түріндегі нәтижесі

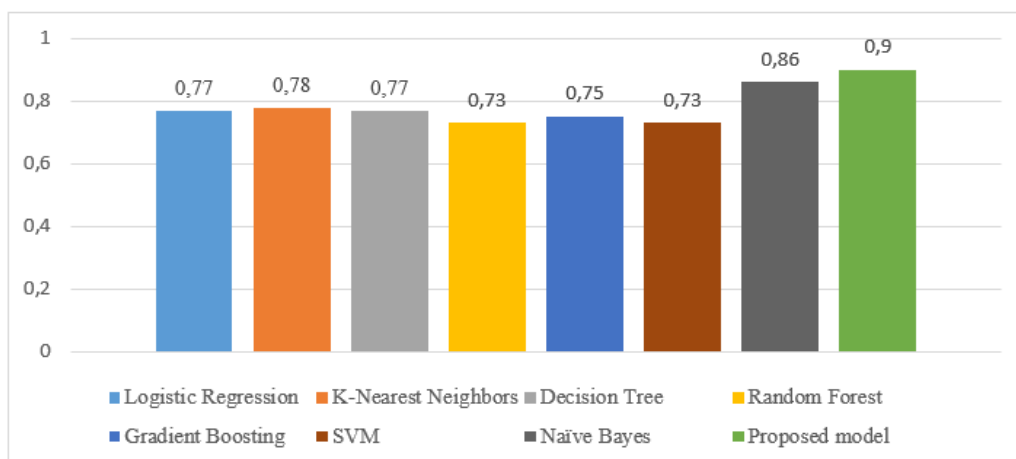
AUC-ROC шамасы төмендегі сурет 4.22-де көрсетілген.



Сурет 4.22 – Градиентті бустинг әдісі арқылы жіктеудің AUC-ROC мәні Қорытындылай келе айтатын болсақ, берілген тарауда бірнеше машиналық оқыту әдістері көмегімен мәтінді экстремистік және бейтарап санаттарға жіктеу есебі шешілді. Эксперимент нәтижесі келесі кесте 4.8-де, сурет 4.23-те келтірілген:

Кесте 4.8 – Машиналық оқыту әдістері көмегімен мәтінді экстремистік және бейтарап санаттарға жіктеу нәтижесі

Машиналық/терең оқыту әдісі	Accuracy	F1-Score	AUC-ROC
Logistic regression	0.77	0.63	0.86
k-nearest neighbors	0.78	0.66	0.5
Decision Tree	0.77	0.64	0.73
Random Forest	0.73	0.52	0.68
Gradient Boosting	0.75	0.58	0.71
SVM	0.73	0.53	0.68
Naïve Bayes	0.86	0.84	0.86
Ұсынылатын модель (TF-IDF_bigram_LSTM)	0.9	0.88	0.89



Сурет 4.23 – Машиналық оқыту әдістері көмегімен мәтінді экстремистік және бейтарап санаттарға жіктеу нәтижесі

Қорытындылай келе, ұсынылатын модель машиналық оқыту әдістерімен салыстырғанда барлық бағалау параметрлері бойынша жоғары екендігін байқаймыз және нәтижесінде ұсынылатын модель веб-ресурстардағы қазақ тіліндегі экстремистік мәтіндерді жоғары дәлдікпен анықтай алады деген қорытындыға келеміз.

## ҚОРЫТЫНДЫ

Берілген диссертациялық жұмыста веб-ресурстардағы қазақ тіліндегі экстремистік мәтіндерді анықтаудың семантикалық талдау моделін құруға қатысты жұмыстар жүргізілді және келесідей нәтижелерге қол жеткізілді:

1) Алғаш рет қазақ тіліндегі экстремистік мәтіндерді анықтау үшін машиналық оқыту әдістерін оқытуға және тестілеуге арналған қазақ тіліндегі экстремистік мәтіндер корпусы құрылды;

2) Алғаш рет қазақ тілінің ерекшеліктерін ескере отырып, LSTM желісінің сөзді ендіру қабатына алдын ала стемминг алгоритмі орындалған биграммдарға TF-IDF әдісін қолданумен ерекшеленетін және экстремистік мәтіндерді анықтау дәлдігін арттыратын семантикалық талдау моделі құрастырылды.

3) Белгілер жиынтығын қалыптастырудың сөзді ендіру әдістері мен n-граммдарды комбинациялауға негізделетін және экстремистік мәтіндерді жіктеудің сапасын арттыратын әдіс құрастырылды.

4) Алғаш рет қазақ тілінде экстремистік түйінді сөздердің тізімі құрылды;

5) Құрастырылған модель мен әдістер нәтижесінде қазақ тіліндегі экстремистік мәтіндерді анықтауға арналған бағдарламалық жабдықтама құрылды.

Бұл зерттеудің жаңалығы қазақ тіліндегі экстремистік мәтіндерді анықтау үшін терең нейрондық желі моделін жасау болып табылады. Алдын ала стемминг алгоритмі орындалған биграммдарға TF-IDF әдісін қолдану негізінде терең нейрондық желі моделі құрастырылды және нәтижелер қазақ тіліндегі экстремистік бағыттағы мәтіндерді анықтау міндеті үшін ең жоғары дәлдікпен классикалық машиналық оқыту әдістерімен салыстырғанда экстремистік мәтіндерді анықтауда ұсынылған модельдің тиімділігін көрсетеді. Осылайша, экстремистік бағыттағы қазақ тіліндегі мәтіндерді табу міндеті үшін биграммдар негізіндегі терең нейрондық желі жоғары өнімділік беріп, өзінің тиімділігін көрсетті.

Диссертациялық жұмыстың теориялық маңыздылығы экстремистік іс-әрекеттер мен ұйымдарды анықтау әдістері мен алгоритмдері саласындағы білім жиынтығына негізделген. Алынған іргелі нәтижелерді әлемдік ғылыми қауымдастық пайдалана алады.

Әдіс, авторлық куәлік түріндегі қолданбалы нәтижелерді ақпараттық қауіпсіздікті, сыни инфрақұрылымды қамтамасыз ету, интернет-экстремизммен күрес жөніндегі уәкілетті органдар пайдалануы мүмкін.

## ПАЙДАЛАНЫЛҒАН ӘДЕБИЕТТЕР ТІЗІМІ

- 1 Gaikwad M., Ahirrao S., Phansalkar S., Kotecha K. Online Extremism Detection: A Systematic Literature Review with Emphasis on Datasets, Classification Techniques, Validation Methods, and Tools // IEEE Access. – 2021. – Vol.9. – P. 48364 – 48404.
- 2 Sánchez-Rebollo C., Puente C., Palacios R., Piriz C., Fuentes J.P., Jarauta J. Detection of Jihadism in Social Networks Using Big Data Techniques Supported by Graphs and Fuzzy Clustering // Complexity. – 2019. – Vol.2019. – P.1 – 13.
- 3 Қазақстан Республикасының №31 Заңы. Экстремизмге қарсы іс-қимыл туралы: 2005 жылдың 18 ақпанында бекітілген. <https://adilet.zan.kz/kaz/docs/Z050000031>
- 4 Қазақстан Республикасы Үкіметінің №124 қаулысы. Қазақстан Республикасында діни экстремизм мен терроризмге қарсы іс-қимыл жөніндегі 2018 – 2022 жылдарға арналған мемлекеттік бағдарламаны бекіту туралы: 2018 жылдың 15 наурызында бекітілген. <https://adilet.zan.kz/kaz/archive/docs/P1800000124/15.03.2018>
- 5 ҚР тыйым салынған шетелдік ұйымдардың тізімі. [https://egov.kz/cms/kk/articles/religion/zaprewennye\\_organizacii?mobile=no](https://egov.kz/cms/kk/articles/religion/zaprewennye_organizacii?mobile=no). 30.03.2021.
- 6 Ақпараттық қауіпсіздік орталығын құруымыз керек – Тоқаев. <https://egemen.kz/article/255352-aqparattyq-qauipsizdik-ortalyghyn-quruymyz-kerek-toqayev>. 06.03.2022.
- 7 Сирия мен Иракта қанша қазақстандық қамауда отыр? <https://ult.kz/post/siriya-men-irakta-kansha-kazakstandyk-kamauda-otyr>. 05.09.2020.
- 8 СІМ "Жусан" операциясы аясында Сириядан тағы 12 азаматты эвакуациялағанын хабарлады. <https://www.azattyq.org/a/31085416.html>. 05.03.2021.
- 9 Addressing the abuse of tech to spread terrorist and extremist content. [https://blog.twitter.com/en\\_us/topics/company/2019/addressing-the-abuse-of-tech-to-spread-terrorist-and-extremist-c](https://blog.twitter.com/en_us/topics/company/2019/addressing-the-abuse-of-tech-to-spread-terrorist-and-extremist-c). 20.06.2021.
- 10 IBM Watson is AI for smarter business. <https://www.ibm.com/watson>. 20.02.2022.
- 11 IT-соавтор «Платона» создает систему мониторинга соцсетей и предсказания угроз. <https://www.vedomosti.ru/technology/articles/2016/06/17/645694-it-soavtor-platona-sozdaet-sistemu-monitoringa-sotssetei-predskazaniya-ugroz>. 18.09.2021.
- 12 We are the central authority for information technology in the security sector. [https://www.zitis.bund.de/DE/Home/home\\_node.html](https://www.zitis.bund.de/DE/Home/home_node.html). 05.04.2022.
- 13 Sharif W., Mumtaz S., Shafiq Z., Riaz O., Ali T., Husnain M., Choi G.S. An Empirical Approach for Extreme Behavior Identification through Tweets Using Machine Learning // Applied Sciences. – 2019. – Vol.9, №18.
- 14 Fatima F., Nurse J. R. C., Goldsmith M. #ISIS vs #ActionCountersTerrorism: A Computational Analysis of Extremist and Counter-extremist Twitter Narratives // 2020 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW). – Genoa: IEEE, 2020. – P. 438-447.

- 15 Alizadeh M., Weber I., Cioffi-Revilla C., Fortunato S., Macy M. Psychology and morality of political extremists: evidence from Twitter language analysis of alt-right and Antifa // EPJ Data Science. – 2019. – Vol.8, №17.
- 16 Benigni, M., Joseph, K. & Carley, K.M. Mining online communities to inform strategic messaging: practical methods to identify community-level insights // Computational and Mathematical Organization Theory. – 2018. – Vol.24, №2. – P. 224–242.
- 17 Ashcroft M., Fisher A., Kaati L., Omer E., Prucha N. Detecting Jihadist Messages on Twitter // 2015 European Intelligence and Security Informatics Conference. Manchester: IEEE, 2015. –P.161-164.
- 18 Omer E. Using machine learning to identify jihadist messages on Twitter. Master thesis: 2015. – Uppsala University, Department of Information technology.
- 19 Ahmad S., Asghar M.Z., Alotaibi F.M., Awan I. Detection and classification of social media-based extremist affiliations using sentiment analysis techniques // Human-centric Computing and Information Sciences. – 2019. –Vol.9, №24. – P. 1 – 23.
- 20 Mitts T. Countering Violent Extremism and Radical Rhetoric // International Organization. – 2022. –Vol.76, №1. – P.251 – 272.
- 21 Gaikwad M., Ahirrao S., Phansalkar S., Kotecha K. 57222563864; Multi-ideology isis/jihadist white supremacist (Miws) dataset for multi-class extremism text classification // Data. – 2021. – Vol.6, №11.
- 22 Aryuni M., Miranda E., Fernando Y., Kibtiah T.M. An early warning detection system of terrorism in indonesia from twitter contents using naïve bayes Algorithm // Proceedings of 2020 International Conference on Information Management and Technology, ICIMTech 2020. – Bandung: IEEE, 2020. –P. 555 – 559.
- 23 Fernandez M., Asif M., Alani H. Understanding the roots of radicalisation on twitter // 10th ACM Conference on Web Science. – Amsterdam: Association for Computing Machinery, 2018. – P.1 – 10.
- 24 van de Weert A., Eijkman Q. Reconsidering Early Detection in Countering Radicalization by Local Frontline Professionals // Terrorism and Political Violence. – 2021. – Vol.33, №2. – P.397 – 408.
- 25 Wei Y., Singh L. Detecting Users Who Share Extremist Content on Twitter. Surveillance in Action. Advanced Sciences and Technologies for Security Applications. – Springer, 2018. – P.351 – 368.
- 26 Fraiwan M. Identification of markers and artificial intelligence-based classification of radical Twitter data // Applied Computing and Informatics. – 2022. – Vol.18, №1.
- 27 López-Sánchez D., Revuelta J., de la Prieta F., Corchado J.M. Towards the automatic identification and monitoring of radicalization activities in twitter // 13th International Conference KMO 2018. – Zilina: Cham Springer, 2018. – P.589 – 599.
- 28 Benigni M., Carley K.M. From tweets to intelligence: Understanding the Islamic Jihad supporting community on twitter // Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics). – Springer International Publishing, 2016. –P.346 – 355.

- 29 Abd-Elaal A.I.A., Badr A.Z., Mahdi H.M.K. Detecting Violent Radical Accounts on Twitter // *International Journal of Advanced Computer Science and Applications (IJACSA)*. – 2020. –Vol.11, №8. –P.516 – 522.
- 30 Mashechkin I., Petrovskiy M., Tsarev D, Chikunov M. Machine Learning Methods for Detecting and Monitoring Extremist Information on the Internet // *Programming and Computer Software*. – 2019. –Vol.45. –P.99 – 115.
- 31 Kotzé E., Senekal B.A., Daelemans W. Automatic classification of social media reports on violent incidents in South Africa using machine learning // *South African Journal of Science*. – 2020. –Vol.116, №3. –P.1 – 8.
- 32 Soliman G.M.A., Abou-El-Enien T.H.M. Terrorism prediction using artificial neural network // *Revue d'Intelligence Artificielle*. – 2019. –Vol.33, №2. –P.81 – 87.
- 33 Al-Zewairi M., Naymat G. Spotting the Islamist Radical within: Religious Extremists Profiling in the United State // *Procedia Computer Science*. – 2017. –Vol.113. – P.162 – 169.
- 34 Ashraf N., Rafiq A., Butt S., Shehzad S.M.F., Sidorov G., Gelbukh A. YouTube based religious hate speech and extremism detection dataset with machine learning baselines // *Journal of Intelligent and Fuzzy Systems*. – 2022. –Vol.42, №5. –P. 4769-4777.
- 35 Ferrara E., Wang W., Varol O., Flammini A., Galstyan A. Predicting Online Extremism, Content Adopters, and Interaction Reciprocity // *International Conference on Social Informatics*. – Bellevue: Springer, 2016. – P.22 – 39.
- 36 Clemmow C., Schumann S., Salman N., Paul G. The Base Rate Study: Developing Base Rates for Risk Factors and Indicators for Engagement in Violent Extremism // *Journal of Forensic Sciences*. – 2020. –Vol.65, №3. –P.865 – 881.
- 37 Pelzer R. Policing of Terrorism Using Data from Social Media // *European Journal for Security Research*. – 2018. –Vol.3. –P.163 – 179.
- 38 Bouchard M, Davies G., Frank R., Wu E., Joffres K. *The Social Structure of Extremist Websites* – Toronto: Toronto University Press, 2020. –P.167 – 189.
- 39 Al-Khoury D. Radicalisation: old and new a comparative analysis of the Red Brigades and the Islamic State // *Quality & Quantity*. – 2020. –Vol.54, №3. –P. 867 – 885.
- 40 Chen H., Thoms S., Fu T. Cyber extremism in Web 2.0: An exploratory study of international Jihadist groups // *2008 IEEE International Conference on Intelligence and Security Informatics*. –Taipei: IEEE, 2008. –P.98 – 103.
- 41 Scanlon J., Gerber M. Automatic detection of cyber-recruitment by violent extremists // *Security Informatics*. – 2014. –Vol.3, №5. –P.1 – 10.
- 42 Aldera S., Emam A., Al-Qurishi M., Alrubaian M., Alothaim A. Online Extremism Detection in Textual Content: A Systematic Literature Review // *IEEE Access*. – 2021. –Vol.9. –P. 42384 – 42396.
- 43 Mouhssine E., Khalid C. Social Big Data Mining Framework for Extremist Content Detection in Social Networks // *International Symposium on Advanced Electrical and Communication Technologies, ISAECT 2018*. – Rabat-Kenitra: IEEE, 2018. –P.1 – 5.
- 44 Araque O., Iglesias C. A. An Approach for Radicalization Detection Based on Emotion Signals and Semantic Similarity // *IEEE Access*. – 2020. –Vol.8. –P.17877 – 17891.



- 45 Hashemi M., Hall M. Visualization, Feature Selection, Machine Learning: Identifying The Responsible Group for Extreme Acts of Violence // IEEE Access. – 2018. – Vol.6. –P.70164–70171.
- 46 Ryan S., Windisch S., Simi P. Former Extremists in Radicalization and Counter-Radicalization Research. Radicalization and Counter-Radicalization. – Bingley, UK: Emerald Publishing Limited, 2020. –P.209 – 224.
- 47 Koehler D. Recent Trends in German Right-Wing Violence and Terrorism: What Are the Contextual Factors behind ‘Hive Terrorism’? // Perspectives on Terrorism. – 2018. – Vol.12, №6. –P.72–88.
- 48 Deviatkin D., Smirnov I., Solovyev F., Suvorova M., Chepovskiy A. Extremist Text Detection In Social Web // Multi Conference on Computer Science and Information Systems, MCCSIS 2019. –Porto, 2019. –P.344–350.
- 49 Ананьева М.И., Девяткин Д.А., Кобозева М.В., Смирнов И.В. Лингвостатистический анализ текстов экстремистской направленности // Ситуационные центры и информационно-аналитические системы класса 4i для задачи мониторинга и безопасности (SCVRT2015-16). Том 1, 2016. –С.210–213.
- 50 Ананьева М. И., Кобозева М. В., Соловьев Ф. Н., Поляков И. В., Чеповский А. М. О проблеме выявления экстремистской направленности в текстах // Вестник Новосибирского государственного университета. Серия: Информационные технологии. – 2016. Том 14, № 4. –С.5–13.
- 51 Devyatkin D., Smirnov I., Ananyeva M., Kobozeva M., Chepovskiy A., Solovyev F. Exploring linguistic features for extremist texts detection (on the material of Russian-speaking illegal texts) // 2017 IEEE International Conference on Intelligence and Security Informatics (ISI). – Beijing: 2017. –P.188–190.
- 52 Бердникова Т.В. Определение адресованности побуждения в экстремистских материалах (на примерах из интернета) // Теория и практика судебной экспертизы. –2019. Том 14, №3. –С.34–39.
- 53 Abbasi A., Chen H. Applying authorship analysis to extremist-group Web forum messages // IEEE Intelligent Systems. –2005. –Vol.20, №5. –P.67-75.
- 54 AlGhamdi M.A., Khan M.A. Intelligent Analysis of Arabic Tweets for Detection of Suspicious Messages // Arabian Journal for Science and Engineering. –2020. –Vol.45, №8. –P.6021–6032 (2020).
- 55 Aldera S., Emam A., Al-Qurishi M., Alrubaian M., Alothaim A Exploratory Data Analysis and Classification of a New Arabic Online Extremism Dataset // IEEE Access. –2021. –Vol.9. –P.161613–161626.
- 56 Alraddadi R.A., Ghembaza M.I.E.-K. Anti-Islamic Arabic Text Categorization using Text Mining and Sentiment Analysis Techniques // International Journal of Advanced Computer Science and Applications. –2021. –Vol.12, №8. –P.776–785.
- 57 Alghofaili H., Almishari M. Countering Terrorism Incitement of Twitter Profiles in Arabic-Context // 21st Saudi Computer Society National Computer Conference, NCC 2018. – Riyadh: IEEE, 2018. –P.224–229.
- 58 Scanlon, J.R., Gerber, M.S. Automatic detection of cyber-recruitment by violent extremists // Security Informatics. –2014. –Vol.3, №5. –P.1–10.

- 59 Cohen K., Johansson F., Lisa K., Clausen M.J. Detecting Linguistic Markers for Radical Violence in Social Media // *Terrorism and Political Violence*. –2014. –Vol.26, №1. –P.246–256.
- 60 Fredrik J., Lisa K., Magnus S. Detecting Linguistic Markers of Violent Extremism in Online Environments. –Hershey: Information Science Reference (an imprint of IGI Global), 2017. –P.374–390.
- 61 Prentice S., Taylor P.J., Rayson P., Hoskins A., O’Loughlin B. Analyzing the semantic content and persuasive composition of extremist media: A case study of texts produced during the Gaza conflict // *Information Systems Frontiers*. – 2011. Vol. 13. –P. 61–73.
- 62 Gelber K. Terrorist-Extremist Speech and Hate Speech: Understanding the Similarities and Differences // *Ethical Theory and Moral Practice*. –2019. –Vol.22. –P.607–622.
- 63 Ul Rehman Z., Abbas S., Khan M.A., Mustafa G., Fayyaz H., Hanif M., Saeed M.A. Understanding the language of ISIS: An empirical approach to detect radical content on twitter using machine learning // *Computers, Materials and Continua*. –2020. –Vol.66, №2. –P.1075–1090.
- 64 Nouh M., Jason Nurse R.C., Goldsmith M. Understanding the radical mind: Identifying signals to detect extremist content on Twitter // *2019 IEEE International Conference on Intelligence and Security Informatics*. –Shenzhen: IEEE, 2019. –P.98–103.
- 65 Rekik A., Jamoussi S., Hamadou A.B. Violent Vocabulary Extraction Methodology: Application to the Radicalism Detection on Social Media. *Lecture Notes in Computer Science* (including subseries *Lecture Notes in Artificial Intelligence* and *Lecture Notes in Bioinformatics*), 11684 LNAI. –Springer, 2019. –P.97–109.
- 66 Kinnvall, C., Capelos, T. The psychology of extremist identification: An introduction [Editorial] // *European Psychologist*. –2021. –Vol.26, №1. –P.1–5.
- 67 Smith L., Wakeford L., Cribbin T., Barnett J., Hou W.K. Detecting psychological change through mobilizing interactions and changes in extremist linguistic style // *Computers in Human Behavior*. –2020. –Vol.108. –P.1–49.
- 68 Bermingham A., Conway M., McInerney L., O’Hare N., Smeaton A. F. Combining Social Network Analysis and Sentiment Analysis to Explore the Potential for Online Radicalisation // *2009 International Conference on Advances in Social Network Analysis and Mining*. –Athens: IEEE, 2009. –P.231–236.
- 69 Scrivens R., Frank R. Sentiment-based Classification of Radical Text on the Web // *2016 European Intelligence and Security Informatics Conference (EISIC)*. – Uppsala: IEEE, 2016. –P.104–107.
- 70 Azizan S.A., Aziz I.A. Terrorism Detection Based on Sentiment Analysis Using Machine Learning // *Journal of Engineering and Applied Sciences*. –2017. –Vol.12, №3, –P.691–698.
- 71 Asif M., Ishtiaq A., Ahmad H., Aljuaid H., Shah J. Sentiment analysis of extremism in social media from textual information // *Telematics and Informatics*. –2020. – Vol.48. –P.1–20.

- 72 Болатбек М.А. Экстремистік мәтіндерді сентимент талдау арқылы анықтау // Международная научная конференция студентов и молодых ученых «Фараби әлемі». – Алматы: Қазақ университеті, 2020. –Б. 75–76.
- 73 Jain A.S., Agarwal S., Agarwal J. A Survey Over Violent Extremist Detection in Social Media Websites // International Journal of Computer Science and Technology, IJCST. –2017. –Vol. 8, №1. –P.71–72.
- 74 Alvares H., Sarkar S., Shakarian P. Detection of Violent Extremists in Social Media // IEEE International Conference on Data Intelligence and Security (ICDIS-19). – USA: IEEE, 2019. –P.43–47.
- 75 Correa D., Sureka A. Solutions to Detect and Analyze Online Radicalization: A Survey // ArXiv. – 2013. –Vol.1301, №4916.
- 76 Arpinar I.B., Kursuncu U., Achilov D. Social media analytics to identify and counter Islamist extremism: Systematic detection, evaluation, and challenging of extremist narratives online // 2016 International Conference on Collaboration Technologies and Systems, CTS 2016. –Orlando:IEEE, 2016. –P.611–612.
- 77 Torregrosa J., Bello Orgaz G., Martínez-Cámara E., Ser J.D., Camacho D. A survey on extremism analysis using natural language processing: definitions, literature review, trends and challenges // Journal of Ambient Intelligence and Humanized Computing. – 2022. –Vol.7, –P.1–37.
- 78 Litvinova T., Litvinova O. Analysis and Detection of a Radical Extremist Discourse Using Stylometric Tools // The 2018 International Conference on Digital Science. – Budva: Springer, 2018. –P.30 – 43.
- 79 Elovici Y., Shapira B., Last M., Zaafrany O., Friedman M., Schneider M., Kandel A. Detection of Access to Terror-Related Web Sites Using an Advanced Terror Detection System (ATDS) // Journal of the Association for Information Science & Technology. –2010. –Vol.61.2010. –P.405–418.
- 80 Mei J., Frank R. Sentiment crawling: Extremist content collection through a sentiment analysis guided web-crawler // 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM). – Paris: IEEE, 2015. – P.1024–1027.
- 81 Zhang Y., Zeng Sh., Fan L., Dang Y., Larson C. A., Chen H. Dark web forums portal: Searching and analyzing jihadist forums // 2009 IEEE International Conference on Intelligence and Security Informatics. –Richardson: IEEE, 2009. –P.71–76.
- 82 Bouchard M., Joffres K., Frank R. Preliminary Analytical Considerations in Designing a Terrorism and Extremism Online Network Extractor. – Springer: Cham, 2014. –P.171–184.
- 83 Scrivens R., Gaudette T., Davies G., Frank R. Searching for Extremist Content Online Using The Dark Crawler and Sentiment Analysis // Methods of Criminology and Criminal Justice Research (Sociology of Crime, Law and Deviance). – 2019. –Vol.24. – P.179–194.
- 84 Akram M., Nasar A., Rehman A., Misuse of charitable giving to finance violent extremism; A futuristic actions study amidst COVID-19 pandemic // Social Sciences & Humanities Open. – 2021. –Vol.4, №1. –P.1–7.

- 85 COVID-19 and Terrorism in the West: Has Radicalization Really Gone Viral? <https://www.justsecurity.org/75064/covid-19-and-terrorism-in-the-west-has-radicalization-really-gone-viral/detect>. 03.01.2022.
- 86 Hate in the time of coronavirus: exploring the impact of the COVID-19 pandemic on violent extremism and terrorism in the West <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7790481/>. 06.04.2022.
- 87 Aktayeva A., Niyazova R., Muradilova G., Makatov Y., Kusainova U. Cognitive Computing Cybersecurity: Social Network Analysis // Communications in Computer and Information Science. – 2020. –Vol.1140. – P. 28–43.
- 88 Bekmanova G., Yelibayeva G., Aubakirova S., Dyussupova N., Sharipbay A., Nyazova R. Methods for Analyzing Polarity of the Kazakh Texts Related to the Terrorist Threats // Computational Science and Its Applications – ICCSA 2019 - 19th International Conference. –Saint Petersburg: Springer, 2019. –P.717–730.
- 89 Мамырбаев О.Ж., Мухсина К.Ж., Хайрова Н.Ф., Колесник А.С. Лингвистические инструменты выявления криминально окрашенной текстовой информации веб-контента // Вестник казахстанско-британского технического университета. –2018. – Том 36 №46. –С.112–117.
- 90 Мамырбаев О.Ж., Хайрова Н.Ф., Мухсина К.Ж. Қазақ тіліндегі мәтіндердегі қылмыстық мәнді коллакцияларды анықтау // М.Тынышбаев атындағы Қазақ көлік және коммуникациялар академиясының хабаршысы. –2019. –Том 3, №110. –Б.170–175.
- 91 Болатбек М.А., Экстремизм түсінігі. Экстремистік мәтіндерді анықтауға арналған белгілер жинағына шолу // Международная научная конференция студентов и молодых ученых «Фараби әлемі». –Алматы: Қазақ университеті, 2020. –Б. 43–44.
- 92 Global Terrorism Database. <https://www.start.umd.edu/gtd/>. 10.03.2022.
- 93 RAND Database of Worldwide Terrorism Incidents. <https://www.rand.org/nsrd/projects/terrorism-incidents.html>. 02.02.2022.
- 94 Байдулла А.М., Мусиралиева Ш.Ж., Болатбек М.А. Экстремистік топтарды анықтау және талдау // Международная научная конференция студентов и молодых ученых «Фараби әлемі». –Алматы: Қазақ университеті, 2021. – Б.74.
- 95 Project Trace. <https://www.interpol.int/Crimes/Terrorism/Counter-terrorism-projects/Project-Trace2>. 03.03.2022.
- 96 Project Sharaka. <https://www.interpol.int/Crimes/Terrorism/Counter-terrorism-projects/Project-Sharaka>. 03.03.2022.
- 97 Project Scorpius. <https://www.interpol.int/Crimes/Terrorism/Counter-terrorism-projects/Project-Scorpius>. 03.03.2022.
- 98 Chen, H. Dark Web: Exploring and Mining the Dark Side of the Web. –Berlin: Springer, Heidelberg, 2011. –P.1–450.
- 99 Mussiraliyeva Sh., Bolatbek M., Omarov B., Bagitova K. Detection Of Extremist Ideation On Social Media Using Machine Learning Techniques // 12th International Conference on Computational Collective Intelligence. – Vietnam, 2020. – P.743–752.
- 100 Mussiraliyeva Sh., Bolatbek M., Omarov B., Medetbek Zh., Baispay G., Ospanov R. On Detecting Online Radicalization and Extremism Using Natural Language

Processing // 21st International Arab Conference on Information Technology (ACIT'2020). – Egypt, 2020. –P.1–5.

101 Python 3.7.0. <https://www.python.org/downloads/release/python-370/>. 03.03.2022.

102 Шәріпбекова С.Е., Мусиралиева Ш.Ж., Болатбек М.А. Қазақ тіліндегі оң қанатты экстремизмді анықтау үшін веб-контентті жинауға арналған бағдарламалық модуль әзірлеу // Международная научная конференция студентов и молодых ученых «Фараби әлемі». –Алматы: Қазақ университеті, 2021. – Б.118.

103 Маден М.Т., Мусиралиева Ш.Ж., Болатбек М.А. Онлайн ортада экстремизмнің лингвистикалық маркерлерін анықтау // Международная научная конференция студентов и молодых ученых «Фараби әлемі». –Алматы: Қазақ университеті, 2021. – Б.101.

104 Мусиралиева Ш.Ж., Омаров Б.С., Болатбек М.А., Жастай Е. Веб-ресурстардағы қазақ тіліндегі экстремисттік сипаттағы мәтіндерді анықтау // Материалы международной научной конференции в области информационных технологий, посвященной 75-летию профессора У.А.Тукеева. – Алматы: Қазақ университеті, 2021. – С. 98-104.

105 Қалиев Ғ. Тіл білімі терминдерінің түсіндірме сөздігі. – Алматы: Сөздік-Словарь, 2005. –Б.439.

106 Қазақ тілі. Энциклопедия. Алматы: Қазақстан Республикасы Білім, мәдениет және денсаулық сақтау министрлігі, Қазақстан даму институты, 1998 жыл, 509 бет.

107 Tukeyev U., Turganbayeva A., Abduali B., Rakhimova D., Amirova D., Karibayeva A. Lexicon-free stemming for Kazakh language information retrieval. // IEEE 12th International Conference on Application of Information and Communication Technologies (AICT).–Almaty, 2018.–P.1–4.

108 Мусиралиева Ш.Ж., Болатбек М.А., Зият Б.М. Стемминг алгоритмі арқылы экстремисттік мәтіндерді жіктеу дәлдігін арттыру // ҚазҰТЗУ хабаршысы. – 2020. №6 (142).–Б.208–215.

109 Мусиралиева Ш.Ж., Болатбек М.А. Әлеуметтік желідегі экстремисттік мәтіндерді жіктеу дәлдігін грамматикалық қателерді анықтау және түзету арқылы арттыру // Международная научно-практическая конференция «Актуальные проблемы информационной безопасности в Казахстане».–Алматы, 2020.–Б.57–61.

110 Болатбек М.А. Создание словаря экстремистских слов для казахского языка // Международная научная конференция студентов и молодых ученых «Фараби әлемі». –Алматы: Қазақ университеті, 2018. –Б.300.

111 Bolatbek M.A., Mussiraliyeva Sh.Zh., Tukeyev U.A. Creating the dataset of keywords for detecting an extremist orientation in web-resources in the Kazakh language // Al-Farabi Kazakh National University, Journal of Mathematics, Mechanics and Computer Science. –2018. –Vol.1, №97. –P.134–142.

112 Sebastiani, F. Machine Learning in Automated Text Categorization // ACM Computing Surveys. – 2002. –Vol.34, №1. –P.1–47.

113 Мусиралиева Ш.Ж., Болатбек М.А. Веб-ресурстардағы экстремисттік мәтіндерді анықтаудың семантикалық үлгілерін құру және зерттеу // Материалы

международной научной конференции студентов и молодых ученых «Фараби әлемі». – Алматы: Қазақ университеті, 2021. – С. 77.

114 Mussiraliyeva Sh., Omarov B., Yoo P., Bolatbek M. Applying Machine Learning Techniques for Religious Extremism Detection on Online User Contents // CMC – Computers, Materials & Continua. – 2021. – Vol. 70, №1. –P. 915–934.

115 Mussiraliyeva Sh., Bolatbek M., Omarov B., Bagitova K., Alimzhanova Zh. Bigram based Deep Neural Network for Extremism Detection in Online User Generated Contents in the Kazakh Language // International Conference on Computational Collective Intelligence. – Greece, 2021. – P.559-570.

116 Swamy M. N., Hanumanthappa M., Jyothi N. M. Indian Language Text Representation and Categorization Using Supervised Learning Algorithm // 2014 International Conference on Intelligent Computing Applications. –Coimbatore, 2014. – P.406–410.

117 TF-IDF. <https://ru.wikipedia.org/wiki/TF-IDF>. 02.02.2022.

118 Jurafsky D., Martin j.H. Description. Speech and Language Processing (2nd Edition). –Pearson, 2009. –P.698.

119 Liu B. Web Data Mining: Exploring Hyperlinks, Contents, and Usage Data. New Jersey: Prentice Hall, 2007. –P.643.

120 van Dam J. K., Zaytsev V. Software Language Identification with Natural Language Classifiers // 2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER). IEEE: –Vol.1. –P.624–628.

121 Pramokchon P., Piamsanga, P. A Feature Score for Classifying Class-Imbalanced Data // In Computer Science and Engineering Conference (ICSEC). – Khon Kaen, 2014. –P.409–414.

122 Deng, L., Yu D. Deep Learning: Methods and Applications. Foundations and Trends in Signal Processing, 2014. –Vol. 7, №3–4. –P.197–387.

123 Bengio Y. Learning Deep Architectures for AI. Foundations and Trends in Machine Learning, 2009. –Vol.2, №1. –P.1–127.

124 Schmidhuber J. Deep Learning in Neural Networks: An Overview. Neural Networks, 2015. –Vol.61.–P.85.

125 Bengio Y., Courville A., Vincent P. Representation Learning: A Review and New Perspectives // IEEE Transactions on Pattern Analysis and Machine Intelligence. – 2013. –Vol.35, №8. –P.1798–1828.

126 Bengio Y., LeCun Y., Hinton G. Deep Learning // Nature. –2015. –Vol. 521. – P.436–444.

127 Гафаров Ф.М. Искусственные нейронные сети и приложения: учебное пособие. – Казань: Издательство Казанского университета, 2018. –С.121.

128 Mussiraliyeva Sh., Omarov B., Bolatbek M., Ospanov R., Vaispay G., Medetbek Zh., Yeltay Zh. Applying Deep Learning for Extremism Detection // International Conference on Advanced Informatics for Computing Research. – Singapore, 2021. – P.597-605.

129 Бурков А. Машинное обучение без лишних слов. –Санкт-Петербург: Питер, 2020. –С.192.

130 Оценка качества в задачах классификации.  
[http://neerc.ifmo.ru/wiki/index.php?title=%D0%9E%D1%86%D0%B5%D0%BD%D0%BA%D0%B0\\_%D0%BA%D0%B0%D1%87%D0%B5%D1%81%D1%82%D0%B2%D0%B0\\_%D0%B2\\_%D0%B7%D0%B0%D0%B4%D0%B0%D1%87%D0%B0%D1%85\\_%D0%BA%D0%BB%D0%B0%D1%81%D1%81%D0%B8%D1%84%D0%B8%D0%BA%D0%B0%D1%86%D0%B8%D0%B8.02.03.2022.](http://neerc.ifmo.ru/wiki/index.php?title=%D0%9E%D1%86%D0%B5%D0%BD%D0%BA%D0%B0_%D0%BA%D0%B0%D1%87%D0%B5%D1%81%D1%82%D0%B2%D0%B0_%D0%B2_%D0%B7%D0%B0%D0%B4%D0%B0%D1%87%D0%B0%D1%85_%D0%BA%D0%BB%D0%B0%D1%81%D1%81%D0%B8%D1%84%D0%B8%D0%BA%D0%B0%D1%86%D0%B8%D0%B8.02.03.2022.)

131 Ынтықбай Б.Н., Мусиралиева Ш.Ж., Болатбек М.А. Элеуметтік желілердегі қауіпсіздік пен конфиденциалдықты машиналық оқыту тәсілдерін қолдану арқылы талдау // Материалы Международной научной конференции студентов и молодых ученых «Фараби әлемі». – Алматы: Қазақ университеті, 2021. – С.119.

132 Болатбек М.А., Мусиралиева Ш.Ж. Экстремистік мәтіндерді машиналық оқыту әдістері арқылы анықтау // ҚазҰТЗУ хабаршысы. – 2018. №6 (130). –Б.300–304.

133 Шалабаев К., Әліпбай К., Болатбек М., Мусиралиева Ш. Вконтакте элеуметтік желісіндегі экстремистік мәтіндерді анықтау және жіктеу // // ҚазҰТЗУ хабаршысы. – 2019. №5 (135). –Б.80–86.

134 Машинное обучение.  
[https://ru.wikipedia.org/wiki/%D0%9C%D0%B0%D1%88%D0%B8%D0%BD%D0%BD%D0%BE%D0%B5\\_%D0%BE%D0%B1%D1%83%D1%87%D0%B5%D0%BD%D0%B8%D0%B5.17.01.2022.](https://ru.wikipedia.org/wiki/%D0%9C%D0%B0%D1%88%D0%B8%D0%BD%D0%BD%D0%BE%D0%B5_%D0%BE%D0%B1%D1%83%D1%87%D0%B5%D0%BD%D0%B8%D0%B5.17.01.2022.)

135 Машинное обучение.  
[http://www.machinelearning.ru/wiki/index.php?title=%D0%9C%D0%B0%D1%88%D0%B8%D0%BD%D0%BD%D0%BE%D0%B5\\_%D0%BE%D0%B1%D1%83%D1%87%D0%B5%D0%BD%D0%B8%D0%B5.20.01.2022.](http://www.machinelearning.ru/wiki/index.php?title=%D0%9C%D0%B0%D1%88%D0%B8%D0%BD%D0%BD%D0%BE%D0%B5_%D0%BE%D0%B1%D1%83%D1%87%D0%B5%D0%BD%D0%B8%D0%B5.20.01.2022.)

136 Naive Bayes classifier. [https://en.wikipedia.org/wiki/Naive\\_Bayes\\_classifier.04.02.2022.](https://en.wikipedia.org/wiki/Naive_Bayes_classifier.04.02.2022.)

137 Градиентный бустинг — просто о сложном.<https://neurohive.io/ru/osnovy-data-science/gradientyj-busting/#:~:text=%D0%93%D1%80%D0%B0%D0%B4%D0%B8%D0%B5%D0%BD%D1%82%D0%BD%D1%8B%D0%B9%20%D0%B1%D1%83%D1%81%D1%82%D0%B8%D0%BD%D0%B3%20E2%80%94%20%D1%8D%D1%82%D0%BE%20%D1%82%D0%B5%D1%85%D0%BD%D0%B8%D0%BA%D0%B0%20%D0%BC%D0%B0%D1%88%D0%B8%D0%BD%D0%BD%D0%BE%D0%B3%D0%BE,%D0%BF%D1%80%D0%B5%D0%B4%D1%81%D0%BA%D0%B0%D0%B7%D1%8B%D0%B2%D0%B0%D1%8E%D1%89%D0%B8%D1%85%20%D0%BC%D0%BE%D0%B4%D0%B5%D0%BB%D0%B5%D0%B9%2C%20%D0%BE%D0%B1%D1%8B%D1%87%D0%BD%D0%BE%20%D0%B4%D0%B5%D1%80%D0%B5%D0%B2%D1%8C%D0%B5%D0%B2%20%D1%80%D0%B5%D1%88%D0%B5%D0%BD%D0%B8%D0%B9.18.02.2022.>

## ҚОСЫМША А

Қазақстан Республикасының аумағында таратуға тыйым салынған топтар

- 1 islam1allah1
- 2 islam\_0
- 3 uhibbyl.islam
- 4 islamk4z
- 5 dnevnikuhti
- 6 medinaschool
- 7 umma.muslim
- 8 malabis\_man
- 9 muslim\_humor
- 10 islam\_today\_ru
- 11 el.islam
- 12 schastya\_v\_islame
- 13 islam\_world
- 14 in\_islams
- 15 hikmaislamicshop
- 16 islam.muslims
- 17 zcalm
- 18 ilovemuhammadiloveislam
- 19 risalatu
- 20 intimislam
- 21 islam\_alhamdulillah
- 22 my.religion.islam
- 23 mydinislam
- 24 muslimsislam
- 25 islamideology
- 26 islamy
- 27 m.dinislam
- 28 koran.sunna
- 29 islam\_obshina
- 30 prayness
- 31 one.ummah
- 32 muslim\_records
- 33 yes\_islam
- 34 taqwaonline
- 35 cumalar\_islam
- 36 istina.v.islame
- 37 v\_islame
- 38 islamnapominanie
- 39 muslim\_record
- 40 raballaha1
- 41 dostovernye.hadisy



- 42 annisa\_today\_ru
- 43 muslimir
- 44 xijra
- 45 uhybbil.islam
- 46 darulkufaridarulislamakhkami
- 47 Өлім-бар умытпа!!
- 48 Өлім Мен Өмірдің Арасы
- 49 Өмір мен өлім
- 50 Аллаһты Ұлықтау
- 51 Құран\_Хадис\_Ислам.
- 52 Құран,хадис және мәтіндер жаттау
- 53 Ислам - әлемдерге мейірім
- 54 Джихад ﷲ
- 55 Аль Джихад
- 56 КНУ Джихад
- 57 КНУ Джихад
- 58 Rasta джихад
- 59 الجهاد Джихад
- 60 ДЖИХАД!
- 61 джихад
- 62 Джихад-OnLine
- 63 Такфир (резерв)
- 64 Шахид
- 65 СУБХАН-АЛЛАХ
- 66 АНТИ-ХАРИДЖИЗМ. ИГИЛ (ИГИШ), аль-Каида
- 67 ДЖИХАД!!!!
- 68 Аль-Каида
- 69 Subhan'Allah.
- 70 СубхьянАллахl
- 71 ALLAHU АКВАР
- 72 АЛЛАХ АКБАР
- 73 Islam
- 74 ИСЛАМ
- 75 Islamic Dinn
- 76 Шейх уль Ислам Ибн Теймия

**ҚОСЫМША Ә**  
Діни мазмұндағы топтар

- 1 abu\_ali\_al\_hanafi
- 2 al\_quran\_kz
- 3 alhamdulilax.musilmanmin
- 4 aliislam0014
- 5 allah\_5
- 6 allahuakbar\_1
- 7 allatagala
- 8 allhamdulillahmusilman
- 9 arman\_ustaz
- 10 atyraumuslims
- 11 bakhyt\_islamda
- 12 club\_musilman
- 13 clubmuslim121212
- 14 cvetok40717
- 15 darulahnaf\_com
- 16 diary\_muslimah\_kz
- 17 dinislam.kz\_toby пусто
- 18 dinislamdinim
- 19 erkanat77777
- 20 hadis110786788
- 21 hadister
- 22 23 hanafi\_kz
- 23 hanafi\_muslims
- 24 hanafi.aktau
- 25 imankz12
- 26 imannyry
- 27 immusulman1
- 28 islaaaaaam
- 29 islam\_1223
- 30 islam\_2017
- 31 islam\_oraza
- 32 islam.offical
- 33 islam2282774
- 34 islamandmuslim
- 35 islamnegizi
- 36 islamsunetti
- 37 issledovanie\_s
- 38 jannatka\_bir\_kadam
- 39 jannatkabirgeinshaallah
- 40 jeckiechankz01
- 41 karagur00

42 kausar\_kz  
43 kaz\_kaz1  
44 kazakh.muslim  
45 kazmuslims  
46 kuran\_jane\_sunnet  
47 kuran\_tadjuid  
48 kuran.hadis  
49 kurannury  
50 kz\_\_sabr  
51 kz\_musilman\_kz  
52 love.south  
53 m\_bauyrlar  
54 m\_kundeligi  
55 m.b.s.kz1111  
56 makatmeshiti  
57 medresekz  
58 men\_musilmanmyn  
59 men.musylmanmyn  
60 menmuslim  
61 mirosulanbiakz  
62 mkk\_kz  
63 muhammad\_kz  
64 musilman\_bauirlar  
65 musilman\_bauirlarim  
66 musilman\_kyndeligi  
67 musilman.hanshaimi  
68 musilman.korgani  
69 musilman.otbasilar.ushin  
70 musilmanbauirlar  
71 musilmanbolubakhyt  
72 musilmandar\_mekeni  
73 musilmandar\_mekenu7117  
74 musilmannin\_sozderi  
75 muslim\_kazakh  
76 muslim362  
77 muslim118104283  
78 muslim12345678910  
79 musliman\_kundelygi  
80 muslimdiarykz  
81 muslimdnevnik  
82 muslimjamagat  
83 muslimpar  
84 muslman.mekeni.jumak  
85 musulman\_alemi\_kz

86 musulman\_allhamdulillah  
87 musulman\_kundeligi  
88 musulman\_kz99  
89 musulmandar\_ordasi  
90 musylman\_kz  
91 musylman\_paryzy  
92 musylman\_sipaty  
93 musylmandar\_parakshasi  
94 musylmankaz  
95 musylmanmyn\_02  
96 muxammedumati  
97 mysilman\_ainasi  
98 namaz\_xanafi  
99 oraza2017  
100 paigambar\_mektebi  
101 paigambar\_omiri  
102 ppageofamuslim  
103 pub.paigambar  
104 quran\_hadith  
105 sizder\_uwin  
106 story\_kaz  
107 surak\_jauap\_hanafi\_mazhabi  
108 tatti\_sezim  
109 wearemuslima  
110 yaasinn  
111 yqylas\_musylman  
112 zhanm\_zhanm  
113 alhamdulillah\_musylmanbyz  
114 musulman.adebi  
115 musulman.kundeligi  
116 1slam\_d1n1  
117 abu\_hanifa000  
118 adamzattanu  
119 ahlusunna\_99  
120 alhalimu  
121 asildin  
122 asyl\_din  
123 baxonbbb  
124 beineu\_kausar  
125 clubislamzhumak  
126 clubmuxambet  
127 din\_islam\_dingegim  
128 din\_jane\_dastur  
129 din\_nasihat

130 din.isla  
131 din.islam\_dingegim  
132 din.islam.dingegim  
133 din.zholy  
134 dinangimeleri  
135 dindastyr  
136 dini\_angymeler  
137 dini\_zhazylymdar  
138 dini\_gibratnama  
139 dinim.islam  
140 dintiregi  
141 dinturaly  
142 dinykitaptar  
143 enu.dt\_kz  
144 f4zhigitteri  
145 iman\_zhuregimde  
146 imandyadam  
147 imannegizderi  
148 immusulman  
149 inshaalla\_sabr  
150 islam\_\_dini\_\_0001  
151 islam\_\_kz  
152 islam\_abu\_hanifa  
153 islam\_akikat\_dini  
154 islam\_akikat  
155 islam\_dini\_jaiynda  
156 islam\_dini1  
157 islam\_dini7  
158 islam\_dinim  
159 islam\_haqq\_din  
160 islam\_iman\_kz  
161 islam\_kazakh  
162 islam\_kzt  
163 islam\_the\_best2  
164 islam\_the  
165 islam.dini.alemi  
166 islam.dini00  
167 islam.islam.dini  
168 islam.xadis  
169 islam120603  
170 islam97925649  
171 islam134849541  
172 islama\_kz  
173 islamdini2015

174 islamdinikzz  
175 islamdinim  
176 islamdinimiz  
177 islamgrupp  
178 islamjumak0  
179 islamkuran  
180 islamkz01  
181 islamkz1400  
182 islammadenieti  
183 islamqazaqstan  
184 jan\_top  
185 jannatka\_asigindar  
186 k\_qasqyr  
187 menin\_dinim\_islam\_kz  
188 menin\_dinimislam  
189 musulman\_balasi\_kz  
190 musulman\_dini  
191 musylman\_kz01  
192 namaz\_dinnin\_tiregy  
193 namaz161566269  
194 namazdintiregi  
195 okomeshit  
196 qissa\_angimeler  
197 sabr\_\_\_kzz  
198 sovkaz  
199 sura114\_hadis6666  
200 taza\_din\_islam  
201 tuka00  
202 allahulyq  
203 almaty\_din  
204 asyldinislamkz  
205 atadinimislam  
206 barlygy\_islam\_dini\_jainda  
207 birikken.odak  
208 d\_islam  
209 darul\_kitab  
210 dasturlydinislam  
211 din.alemi  
212 din.islam.dingegum  
213 din.religion.semey  
214 din\_islam\_1  
215 dini.angimeler  
216 dini\_angimeler147684564  
217 diniangimeler

218 dinim\_islam  
219 dinimislam  
220 dinisauat  
221 dinislam231217  
222 dinkgu\_oral  
223 dintanym  
224 dintube  
225 en.taza.islam.dini  
226 gibryatnama\_dini\_angimeler  
227 idkontrol00  
228 iman\_islam0  
229 islam\_dini\_2019  
230 islam\_dini\_alla\_zholyinda  
231 islam\_dini\_taza\_din  
232 islam\_dini.khak  
233 islam\_the\_best  
234 islam7242  
235 islam011111  
236 islam84550948  
237 islameli2015  
238 Islamikz  
239 Islamjoly  
240 Islamjumak  
241 Islamsunnakz  
242 islamtv2  
243 Jamagatkz  
244 kalash\_kz  
245 kaz.uagiz  
246 kazakhmuslim  
247 kitap\_alemi  
248 kmdb\_vk  
249 menim\_dinim\_islam  
250 menin\_dinim\_islam\_1  
251 muslims010  
252 Musulmanyislam  
253 namaz\_dinnin\_tiregi  
254 namaz.dinnin.tiregi  
255 Namazzaman  
256 ne\_poim  
257 official\_kz  
258 qazaq\_islam  
259 qiziktv77777  
260 qmdb\_dinikitaptar  
261 trapmuz\_kz

262 tur\_dmzo  
263 udralmatinskayaoblast  
264 xktu\_din\_teo  
265 yslamdini\_kz



## ҚОСЫМША Б

### Жалпы лексиканы қамтитын топтар

- 1 Пайдалы Кеңестер pajdaly\_kenester\_kz
- 2 ӘЗІЛ ӘЛЕМІ azil\_qz
- 3 Әзіл әлемі azil.alemi2015
- 4 Әзілдер azil\_\_der
- 5 Әзіл-қалжың kalzhin
- 6 ӘЗІЛ ӘЛЕМІ kazworld
- 7 Жынды әзілдер z\_azilder
- 8 ӘЗІЛ ОРДАСЫ azilo\_kz
- 9 Ең үздік әзілдер enuzdikazilder
- 10 Қазақ жігіттері kazakh\_kz\_jigitter
- 11 Неке Қазақстан nikah\_kazahstan
- 12 Жаңа қазақша әндер порталы!!! newmusic\_kz
- 13 Қазаққа жақпайың... kazakka\_jakpaisyn
- 14 Қазақша есімдер kaz.esimder
- 15 Қазақша анекдоттар anekdottar\_goi
- 16 Қазақстан Вконтакте kz
- 17 Менің ойларым|Қазақша Бот azil\_mekeni
- 18 Қазақ Футболы qazaq\_football
- 19 Тек қана қазақша әзілдер tek.kz\_azil
- 20 Әйелдерге арналған ayelge
- 21 Өз қолыңмен/ Пайдалы кеңестер oz\_kolinmenn
- 22 Фактілер Әлемі kz\_logick
- 23 Психология әлемі aleyu
- 24 Әлем тарихы alem.history
- 25 Футбол Әлемі footballinkz
- 26 Әйелдер әлемі ayelderalemy
- 27 Аспазшылар әлемі! as\_xanakz
- 28 Поэзия әлемі kazpoetry
- 29 Философия Әлемі kzphylosophy
- 30 Әлем Футболы footballinqaz
- 31 Жұлдыздар әлемі starsalem
- 32 Кітаптар әлемі kitaptar
- 33 Әлем Кереметтері kazakh\_alm
- 34 Дәмхана damxana
- 35 Жалғыздық jalgyz\_dyk
- 36 Махаббат|Отбасы|Балалар mahabbat\_balalar
- 37 Отбасы сырлары otbasysyrlary
- 38 DalaFilm | DF dalafilm
- 39 ТОМПИ tompi06021995
- 40 ТОМПИ tompi\_\_tompi
- 41 5 ҚЫЗЫҚТЫ ФАКТІЛЕР 5kizik

42 Мектеп қызықтары mektep\_kyzyktary  
 43 Ұлы сөздер U.S uly\_sozder  
 44 Махаббат,қызық мол жылдар maxabbat\_kz  
 45 Сіз білмейтін қызықтар kzktar  
 46 ҚызықТаймыз kiziqtaimiz\_100k  
 47 Сіз білмейтін қызықтар kiziktar  
 48 Бір миллион қызық оқиға millionokiga  
 49 Сіз білмейтін қызықтар sbkiziktar  
 50 Қызықты Дерктер |KD kz\_facts  
 51 Қызықтың бәрі осында kyzyktynbary  
 52 Әлем қызықтары alemkyzyktari  
 53 Қызықсыз факттер quirdaq  
 54 Бекболат Тілеухан btleukhan  
 55 Ағылшын тілін үйрен|QAZENG qazeng  
 56 Ағылшын тілін үйренейік) publicgreatbl7  
 57 Ағылшын тілі| Күн сайын engkazlanguage  
 58 ҰБТ Академиясы ubt\_akademiyasy  
 59 "QAZENT" –новый формат ЕНТ 2020 ҰБТ|Kazent kazent  
 60 АБАЙ ҚҰНАНБАЙҰЛЫ/АБАУ QUNANBAYULY abay\_qunanbayuly  
 61 Омар Хайям және ұлы философтар omarzhene  
 62 Анекдоттар мен приколдар kz\_version  
 63 Мықты приколдар myktyprikol  
 64 Жынды приколдар станциясы JPS jpskz  
 65 Қош махаббат qoshmahabbat  
 66 Бітпейтін махаббат bitpeityn\_maxabbat  
 67 Сұлу махаббат sulumahabbat  
 68 Махаббат сезімі mahabbat\_sezimi  
 69 Махаббат кілті –жүректе mahabbat\_kilti  
 70 Өмір және Махаббат... kaz\_\_world  
 71 Мөлдір Махаббат hubbi  
 72 Тентек Сезім tenteksezym  
 73 Сезім sezim  
 74 Әйел сезімі ayelsezymy  
 75 Қыз сезімі qiz\_sezimi  
 76 NUR.KZ-жаңалықтар portalnurkz\_kaz  
 77 Push-Жаңалықтар push\_kz  
 78 Жаңалықтар qazaqshanews  
 79 Жаңалықтар KZ qazaq\_news  
 80 Өмірден алынған omirden\_allingan  
 81 Өмір күнделігі omirkundeligy  
 82 Өмірде бәрі мүмкін kzomirde  
 83 Өмірлік кеңестер omir\_kenes  
 84 Өмірде солай omirdesolay  
 85 Өмір жайлы... omirj\_kz

- 86 Жұлдыздардың өмірі juldyzomir
- 87 Өмір lifeomir
- 88 Өмір аялдамасы omirayaldama
- 89 Өмірғой omirlik.sirlasyn
- 90 Өмір omir\_7
- 91 Менің өмірім omiringoi
- 92 Өмір жайлы бірер сөз omirzhaily
- 93 Менің өмірім omir02
- 94 Есіңдеме достым e\_dostym
- 95 Өмірдің ащы шындық shindik\_awi
- 96 Өмірдегі сөздер omirdegi\_sozder
- 97 Өмірлік кеңестер omirlyk\_kenester
- 98 Жұлдызды әлем zh\_alem
- 99 Жұлдыз жорамал|Күн сайын zhuldiz.zhoranal
- 100 Үздік жұлдыз жорамал|Күн сайын juldiz\_joramall

## ҚОСЫМША В

Экстремистік әрекет туралы жаңалықтар мақалаларына сілтеме

- 1 <https://qazaquni.kz/2019/10/12/105733.html>
- 2 <https://qazaquni.kz/2019/09/20/104819.html>
- 3 <https://qazaquni.kz/2019/06/03/100587.html>
- 4 <https://qazaquni.kz/2019/04/06/98048.html>
- 5 <https://qazaquni.kz/2019/01/09/94470.html>
- 6 <https://qazaquni.kz/2018/10/29/91726.html>
- 7 <https://qazaquni.kz/2018/01/30/80985.html>
- 8 <https://qazaquni.kz/2017/08/10/73397.html>
- 9 <https://qazaquni.kz/2017/04/07/67449.html>
- 10 <https://qazaquni.kz/2016/12/25/62454.html>
- 11 <https://qazaquni.kz/2016/12/21/62341.html>
- 12 <https://qazaquni.kz/2016/10/11/59057.html>
- 13 <https://qazaquni.kz/2016/09/26/58512.html>
- 14 <https://qazaquni.kz/2016/09/26/58484.html>
- 15 <https://qazaquni.kz/2016/08/31/57154.html>
- 16 <https://qazaquni.kz/2016/07/22/55667.html>
- 17 <https://24.kz/kz/zha-aly-tar/sayasat/item/363424-m-zhilis-shy-ekstremizmmen-k-res-konventsiasyn-ma-ldady>
- 18 <https://24.kz/kz/zha-aly-tar/o-am/item/349748-sham-elinen-kelgenderdi-shuly-soty-li-de-zhal-asady>
- 19 <https://24.kz/kz/zha-aly-tar/o-am/item/348370-aza-standa-22-dini-a-ymny-taraluyna-git-nasikhatyna-ata-tyjym-salyn-an>
- 20 <https://24.kz/kz/zha-aly-tar/o-am/item/277572-sarapshylar-ekstremizmni-aldyn-aluda-da-bilimni-lesi-zor>
- 21 <https://24.kz/kz/zha-aly-tar/o-am/item/273840-ma-ystau-oblysynda-y-t-zetu-mekemesinde-o-shaulan-an-ajma-ryldy>
- 22 <https://24.kz/kz/zha-aly-tar/o-am/item/247224-sarapshy-zhastardy-ziyandy-ideologiyany-serinen-or-audy-zholy-bilim-beru>
- 23 <https://24.kz/kz/zha-aly-tar/o-i-a/item/238961-alardy-tarat-an-8-azamat-staldy>
- 24 <https://24.kz/kz/zha-aly-tar/o-am/item/227531-almatyda-studentterge-dini-ekstremizmni-aupi-zh-ne-onymen-k-res-zholdary-t-sindirildi>
- 25 <https://24.kz/kz/zha-aly-tar/o-am/item/216913-kadam-shakh-shakhim-aza-standau-anstan-m-selesin-sheshu-shin-ta-y-bir-lken-adam-zhasady>
- 26 <https://24.kz/kz/zha-aly-tar/o-am/item/212185-aza-standa-men-t-rikmenstan-khaly-araly-terrorizmge-arsy-k-reste-riptestikti-k-shejtpek>
- 27 <https://24.kz/kz/zha-aly-tar/o-am/item/210501-astanada-ba-kilderi-men-memlekettik-yzmetkerler-bas-osty>
- 28 <https://24.kz/kz/zha-aly-tar/lemde/item/184204-tmd-zhastaryny-parlamentaraly-assambleyasynnda-terrorizmge-arsy-t-ru-arastyryldy>
- 29 <https://24.kz/kz/zha-aly-tar/o-am/item/177894-ekstremizmdi-nasikhattajtyyn-300-sajt-zhabyldy>

- 30 <https://24.kz/kz/archive/zha-aly-tar/item/150878-ekstremizidi-zhe-uge-bolady>
- 31 <https://24.kz/kz/archive/zha-aly-tar/item/145664-zirbajzhan-diasporasyny-zhastary-kazakhstan-for-peas-oz-alyssyna-osyldy>
- 32 <https://24.kz/kz/archive/zha-aly-tar/item/145420-ibitsa-aralynda-zhastardy-terrorister-ataryna-tart-an-imamdar-staldy>
- 33 <https://24.kz/kz/archive/zha-aly-tar/item/144231-khaly-araly-konferentsiya-bastaldy>
- 34 <https://24.kz/kz/archive/zha-aly-tar/item/143013-nigeriyada-boko-kharam-21-yzdy-bosatty>
- 35 <https://sn.kz/sn-akparat-agyny/63556-almatyda-en-uzdik-leumettik-rolik-avtorlary-marapattaldy>
- 36 <https://sn.kz/sn-zan-zaman/56462-ruk-sat-etilmegen-miting-otkizuge-daiyndalghan-dvk-kozgalysynyn-tort-belsendisi-kamauga-alyndy>
- 37 <https://sn.kz/sn-sayasat/51122-yr-yzstanda-ekstremistik-sipatta-y-36-sajt-b-attaldy>
- 38 <https://sn.kz/b-gingi-s-z/43617-elimizdegi-bajsaldy-salafilardi-bir-mezette-radikaldy-salafilerge-ajnalyp-ketui-m-mkin>
- 39 <https://sn.kz/sn-zan-zaman/40463-a-t-bede-tabli-i-zhama-at-jymyny-7-m-shesi-sottaldy>
- 40 <https://sn.kz/o-am/37484-shymkentte-terrorizm-babymen-sottal-an-adamny-tuystary-shu-shy-ardy>
- 41 <https://sn.kz/o-am/37050-o-o-da-zh-rtty-terrorizmge-gittegen-ta-y-2-adam-staldy>
- 42 <https://sn.kz/o-am/36503-o-o-da-zh-rtty-terakt-zhasau-a-ndep-zh-rngen-sa-aldylar-staldy>
- 43 <https://sn.kz/sn-akparat-agyny/32341-t-zhikstanda-70-zhasta-y-zhej-nemerelerimen-birge-siriya-a-so-ys-a-attanba-bol-an>
- 44 <https://sn.kz/sn-zan-zaman/31086-terrorizmdi-nasikhattady-dep-ajyptal-an-shymkenttik-zhigit-asyl-arnadan-o-yl-an-s-reni-zh-ktep-al-anyn-ajtty>
- 45 <https://sn.kz/bilik/30021-nazarbaev-auipsizdik-ke-esin-zhinap-o-o-da-y-dini-ekstremizm-m-selesin-tal-ylady>
- 46 <https://sn.kz/b-gingi-s-z/29845-atyrauda-sottal-an-salafit-daryn-m-b-rovty-ty-damau-a-sha-yrdy>
- 47 <https://sn.kz/bilik/29745-nazarbaev-pen-m-simov-dini-radikaldar-zhajynda-gimelesti>
- 48 <https://sn.kz/sn-zan-zaman/29169-ltty-auipsizdik-komiteti-o-o-da-bir-top-sa-aldy-azamaty-stady>
- 49 <https://sn.kz/b-gingi-s-z/28835-za-ger-salafittik-a-ym-a-erip-ketken-zhastardy-obaly-e-birinshiden-daryn-basta-an-topty-mojnynda>
- 50 <https://sn.kz/o-am/8539-prokuror-ruslan-k-lekbaev-a-lim-zhazasyn-s-rady>
- 51 <https://sn.kz/o-am/7966-elimizdi-batysynda-la-kestik-zhasama-bol-an-21-adam-staldy>
- 52 <https://sn.kz/sn-akparat-agyny/6768-siriya-birneshe-zharylys-bolyp-119-adam-aza-tapty>

- 53 <https://sn.kz/sn-akparat-agyny/6703-astanada-siriya-a-otbasymen-attanba-bolan-eki-t-r-yn-sottaldy>
- 54 <https://sn.kz/sn-akparat-agyny/6570-ekstremistik-k-z-arasta-y-ta-y-bir-aza-standy-saudiya-koroldiginde-ol-a-t-sti>
- 55 <https://sn.kz/sn-akparat-agyny/6385-islam-memleketi-400-adamdy-t-t-yn-aldy>
- 56 <https://sn.kz/sn-akparat-agyny/6289-1-kaida-italiya-men-ispaniya-a-shabuyl-zhasaudy-zhosparlap-otyr>
- 57 <https://sn.kz/sn-akparat-agyny/6214-sh-yl-saud-arabiyasynda-ekstremistik-k-z-arasta-y-aza-stan-azamattary-staldy>
- 58 <https://sn.kz/sn-akparat-agyny/5977-islam-memleketini-kapitaly-1-5-mlrd-dollardy-rajdy>
- 59 <https://sn.kz/sn-akparat-agyny/5875-islam-memleketi-zhas-sodyrlardy-dayarlajtyn-zh-je-rma>
- 60 <https://sn.kz/o-am/5550-m-zhilis-deputaty-k-rim-m-simovten-la-kestikke-arsy-k-resti-k-shejtudi-s-rady>
- 61 <https://sn.kz/o-am/1806-ekstremizmni-aldyn-aludy-y-ty-negizderi>
- 62 <https://sn.kz/o-am/1804-ara-andyly-medsina-yzmetkerlerine-dini-radikalizmge-arsy-t-ru-zholdary-t-sindirildi>
- 63 <https://sn.kz/sn-zan-zaman/1687-temirtauda-terrorizm-zh-ne-ekstremizmmen-k-res-m-seleleri-tal-ylandy>
- 64 <https://sn.kz/o-am/541-a-yz-adam-zhurnalyny-redaktoryna-ekstremizmdi-atady-degen-zha-a-ajyp-ta-yldy>
- 65 [https://www.inform.kz/kz/ekstremizm-terrorizm-separatizm-dinge-zhat-ugym\\_a3572197](https://www.inform.kz/kz/ekstremizm-terrorizm-separatizm-dinge-zhat-ugym_a3572197)
- 66 <http://prokuror.gov.kz/kaz/baspasoz/makalalar/ekstremizm-men-terrorizmge-karsy-kimyl-kogam-turaktylygynyn-kepili>
- 67 <https://azh.kz/kz/news/view/8318>
- 68 [https://www.inform.kz/kz/zhastar-arasyndagy-dini-ekstremizmnin-kalay-aldyn-aluga-bolady\\_a3572203](https://www.inform.kz/kz/zhastar-arasyndagy-dini-ekstremizmnin-kalay-aldyn-aluga-bolady_a3572203)
- 69 [https://kaz.tengrinews.kz/kazakhstan\\_news/sagyintaev-ekstremizm-men-terrorizmge-karsyi-s-kimiyil-275091/](https://kaz.tengrinews.kz/kazakhstan_news/sagyintaev-ekstremizm-men-terrorizmge-karsyi-s-kimiyil-275091/)
- 70 <https://ustinka.kz/kz/kazakhstan/politics/35278.html>
- 71 <https://sputniknews.kz/incidents/20190602/10353234/Almaty-ekstremizm-usher-adam.html>
- 72 [https://kazakh-tv.kz/kz/view/news\\_kazakhstan/page\\_172092\\_ekstremizm-men-terrorizmge-toskauyl-koyu-kazhet](https://kazakh-tv.kz/kz/view/news_kazakhstan/page_172092_ekstremizm-men-terrorizmge-toskauyl-koyu-kazhet)
- 73 <http://www.nkzu.kz/news/view?id=6593>
- 74 <https://www.ktk.kz/kz/programs/novosti/70984/>
- 75 <https://adyrna.kz/post/20826>
- 76 <https://egemen.kz/article/167748-ekstremizm-men-terrorizmninh-aldyn-aludgayy-talqylandy>
- 77 [https://www.kt.kz/kaz/state/kazakstanda\\_ekstremizm\\_zhane\\_terrorizm\\_kilmisi\\_na\\_karsi\\_kilmistik\\_zhauapkershilik\\_kushejtiledi\\_1153626546.html](https://www.kt.kz/kaz/state/kazakstanda_ekstremizm_zhane_terrorizm_kilmisi_na_karsi_kilmistik_zhauapkershilik_kushejtiledi_1153626546.html)

- 78 <https://www.vuzkunaeva.kz/index.php/kz/main-news-kz/666-otchet-o-provedenii-kruglogo-stola-na-temu-protivostoyanie-religioznomu-ekstremizmu-i-terrorizmu-problemy-i-puti-resheniya-kz>
- 79 [https://mks.gov.kz/press-sluzhba/novosti\\_ministerstva/?cid=0&rid=1703](https://mks.gov.kz/press-sluzhba/novosti_ministerstva/?cid=0&rid=1703)
- 80 <https://khabar.kz/kk/muragat-ajbyn/item/69847-ajbyn-aza-standa-ekstremizm-zh-ne-terrorizmmen-k-resetin-arnajy-zhasa-tar>
- 81 [http://old.baq.kz/kk/news/aimaktik\\_bak\\_muragat/regmedia-43237](http://old.baq.kz/kk/news/aimaktik_bak_muragat/regmedia-43237)
- 82 [https://kaz.tengrinews.kz/kazakhstan\\_news/kazakstanda-biyil-3-terraktnn-jolyi-kesld-ukk-303537/](https://kaz.tengrinews.kz/kazakhstan_news/kazakstanda-biyil-3-terraktnn-jolyi-kesld-ukk-303537/)
- 83 [https://kaz.tengrinews.kz/kazakhstan\\_news/jusan-operatsiyasyi-jalgasadyi-303348/](https://kaz.tengrinews.kz/kazakhstan_news/jusan-operatsiyasyi-jalgasadyi-303348/)
- 84 <https://kaz.tengrinews.kz/crime/turkstanda-lankestkke-undegen-turgyin-sottaldyi-301212/>
- 85 [https://kaz.tengrinews.kz/kazakhstan\\_news/kazakstanda-yakyin-inkar-uyyimyiektremistk-dep-tanyildyi-292726/](https://kaz.tengrinews.kz/kazakhstan_news/kazakstanda-yakyin-inkar-uyyimyiektremistk-dep-tanyildyi-292726/)
- 86 <https://kaz.tengrinews.kz/crime/shyimkentte-terakt-jaylyi-jalghan-habar-bergen-er-adam-297604/>
- 87 <https://egemen.kz/article/200109-saqshy>
- 88 <https://egemen.kz/article/179465-islam-dgane-terrorizm-odaqtas-pa-aldeqarsylas-pa>
- 89 <https://egemen.kz/article/168315-qostanayda-telefon-terroristerininh-qonhyrauynan-kelgen-shyghyn-35-million-tenhg>
- 90 <https://egemen.kz/article/163332-terrorizmmen-kures-%E2%80%93-barshanyh-mindeti>
- 91 <https://egemen.kz/article/165064-bylytyr-zorlyq-zombylyqty-nasikhattaytyn-9-mynhgha-dguyq-sayt-dgumysy-toqtatyldy>
- 92 <https://holanews.kz/view/news/19435>
- 93 <https://holanews.kz/view/news/15040>
- 94 [https://kaz.tengrinews.kz/kazakhstan\\_news/sak-bolyinyizdar-ukk-youtube-jelsnde-video-jariyaladyi-304143/](https://kaz.tengrinews.kz/kazakhstan_news/sak-bolyinyizdar-ukk-youtube-jelsnde-video-jariyaladyi-304143/)
- 95 [https://kaz.tengrinews.kz/kazakhstan\\_news/40-kazakstandyik-shetelde-sodyirlar-kataryinda-jur-ukk-303347/](https://kaz.tengrinews.kz/kazakhstan_news/40-kazakstandyik-shetelde-sodyirlar-kataryinda-jur-ukk-303347/)
- 96 <https://www.ktk.kz/kz/blog/article/2016/10/05/72691/>
- 97 <https://www.ktk.kz/kz/blog/article/2016/07/25/71025/>
- 98 <https://www.ktk.kz/kz/news/video/2015/01/23/56859/>
- 99 <https://www.ktk.kz/kz/news/video/2018/03/13/91918/>
- 100 <https://kaz.nur.kz/1825896-abaev-elimizde-kimderge-ekstremizm-kaupi-tnip-trganyn-mlimdedi.html>
- 101 <https://informburo.kz/kaz/azastan-islamny-ltty-modeln-zhasau-kerek.html>
- 102 <https://www.zakon.kz/4848681-terrorizm-elge-t1257ngen-1179au1110p.html>
- 103 <https://www.zakon.kz/video/v/73690.html>
- 104 <https://www.zakon.kz/video/v/13765.html>
- 105 <https://www.zakon.kz/video/v/3983.html>

- 106 <https://www.zhasalash.kz/article/5657-uqk-qazaqstanda-20-mynhnan-astam-adam-ekstremizmdi-dgaqtaydy>
- 107 <https://www.zhasalash.kz/article/685-alla-men-adam-arasyn-saudalap-dgurgenderdinh-maqsaty-ne>
- 108 <https://www.zhasalash.kz/article/1021-eklektika-ekstremizmi>
- 109 <http://kazaknews.kz/okiga/jekstremizm-men-terrorizm-m-seleleri/>
- 110 <http://kazaknews.kz/tehnologiya/aza-standa-tyjym-salyn-an-18-my-sajt-belgili-boldy/>
- 111 <http://kazaknews.kz/okiga/a-t-be-oblysynda-bylytyrdan-beri-terrorizm-ajybymen-13-adam-t-rmege-amaldy/>
- 112 <http://kazaknews.kz/okiga/k-radikalizmge-arsy-sharalar-zhasa-tajdy-2018-de-3-la-kestik-reketti-aldy-aly-an/>
- 113 <http://kazaknews.kz/kogam/sirija-auma-ynan-kelgen-balalal-dn-synamasynan-tedi/>
- 114 [https://www.azattyq.org/a/kazakhstan\\_students\\_are\\_warned\\_about\\_foreign\\_enemies\\_and\\_terrorist\\_plots/24538460.html](https://www.azattyq.org/a/kazakhstan_students_are_warned_about_foreign_enemies_and_terrorist_plots/24538460.html)
- 115 <https://www.azattyq.org/a/kazakhstan-chinese-kazakhs-who-cross-the-border/30391639.html>
- 116 <https://www.azattyq.org/a/30011733.html>
- 117 <https://www.azattyq.org/a/kazakhstan-aktobe-syria/29998676.html>
- 118 <https://qazaqtimes.com/article/12721>