

**МОНГОЛ УЛСЫН ИХ СУРГУУЛЬ
ХЭРЭГЛЭЭНИЙ ШИНЖЛЭХ УХААН, ИНЖЕНЕРЧЛЭЛИЙН СУРГУУЛЬ
МЭДЭЭЛЭЛ, КОМПЬЮТЕРИЙН УХААНЫ ТЭНХИМ**

Анужингийн Сайнзолбоо

**АЖИЛ ОЛГОГЧДЫН ӨГӨГДЛИЙН АНАЛИЗ
СИСТЕМ ДЭЭР СУУРИЛСАН ЧАТ БОТ
(Chat bot based on sytem analysis of employers' data)**

Мэдээллийн технологи (D061303)
Бакалаврын судалгааны ажил

Улаанбаатар

2022 оны 03 сар

МОНГОЛ УЛСЫН ИХ СУРГУУЛЬ
ХЭРЭГЛЭЭНИЙ ШИНЖЛЭХ УХААН, ИНЖЕНЕРЧЛЭЛИЙН СУРГУУЛЬ
МЭДЭЭЛЭЛ, КОМПЬЮТЕРИЙН УХААНЫ ТЭНХИМ

АЖИЛ ОЛГОГЧДЫН ӨГӨГДЛИЙН АНАЛИЗ СИСТЕМ
ДЭЭР СУУРИЛСАН ЧАТ БОТ

(Chat bot based on sytem analysis of employers' data)

Мэдээллийн технологи (D061303)
Бакалаврын судалгааны ажил

Удирдагч: _____ Б.Хуягбаатар доктор (Ph.D.)
Гүйцэтгэсэн: _____ А.Сайнзолбоо (18B1NUM1762)

Улаанбаатар

2022 оны 03 сар

Зохиогчийн баталгаа

Миний бие Анужингийн Сайнзолбоо “АЖИЛ ОЛГОГЧДЫН ӨГӨГДЛИЙН АНАЛИЗ СИСТЕМ ДЭЭР СУУРИЛСАН ЧАТ БОТ” сэдэвтэй судалгааны ажлыг гүйцэтгэсэн болохыг зарлаж дараах зүйлсийг баталж байна:

- Ажил нь бүхэлдээ эсвэл ихэнхдээ Монгол Улсын Их Сургуулийн зэрэг горилохоор дэвшүүлсэн болно.
- Энэ ажлын аль нэг хэсгийг эсвэл бүхлээр нь ямар нэг их, дээд сургуулийн зэрэг горилохоор оруулж байгаагүй.
- Бусдын хийсэн ажлаас хуулбарлаагүй, ашигласан бол ишлэл, зүүлт хийсэн.
- Ажлыг би өөрөө (хамтарч) хийсэн ба миний хийсэн ажил, үзүүлсэн дэмжлэгийг дипломын ажилд тодорхой тусгасан.
- Ажилд тусалсан бүх эх сурвалжид талархаж байна.

Гарын үсэг: _____

Огноо: _____

	ГАРЧИГ
УДИРТГАЛ	1
БҮЛГҮҮД	2
1. СЭДВИЙН ТАНИЛЦУУЛГА	2
1.1 Оршил	2
1.2 Зорилго	2
1.3 Зорилт	2
1.4 Алсын хараа	3
2. СИСТЕМИЙН СУДАЛГАА	4
2.1 Системийн судалгаа	4
2.2 Ижил төстэй системүүд	7
2.3 Технологийн судалгаа	9
3. СИСТЕМИЙН ШИНЖИЛГЭЭ	16
3.1 Бизнесийн үйл ажиллагааны шинжилгээ	16
3.2 Хэрэглэгч	17
3.3 Функционал шаардлага	17
3.4 Функционал бус шаардлага	18
3.5 Use case диаграмм	19
4. СИСТЕМИЙН ЗОХИОМЖ	20
4.1 Өгөгдлийн сангийн диаграмм	20
4.2 Өгөгдлийн элемент	21
4.3 Өгөгдлийн сангийн холбоосын тайлбар	24
5. ХЭРЭГЖҮҮЛЭЛТ, ҮР ДҮН	25
5.1 Хөгжүүлсэн байдал	25

ДҮГНЭЛТ	39
ДҮГНЭЛТ	39
НОМ ЗҮЙ	40
ХАВСРАЛТ.....	41
А. ҮЕЧИЛСЭН ТӨЛӨВЛӨГӨӨ	41
В. КОДЫН ХЭРЭГЖҮҮЛЭЛТ.....	42
В.1 Өгөгдөл цугуулалт	42
В.2 Өгөгдөл нэгтгэх, цэвэрлэх	48

ЗУРГИЙН ЖАГСААЛТ

2.1	Pizza Hut chat bot	7
2.2	WHO's chat bot	8
2.3	Python лого	9
2.4	BeautifulSoup лого	10
2.5	Өгөгдөл цуглуулалтын жишээний үр дүн	11
2.6	Cosine similarity утгын график	12
2.7	SBERT.net лого	12
2.8	cosine-similarity ашигласан жишээ	13
2.9	Чатботын амьралын мөчлөг	14
3.1	BPMN-1	16
3.2	BPMN-2	17
3.3	Use Case диаграмм	19
4.1	Өгөгдлийн сангийн диаграмм	20
5.1	Үндсэн процесс зураглал	26
5.2	Цуглуулсан өгөгдлийн файлууд	31
5.3	Data set	32
5.4	Өгөгдлийг цэвэрлэж, өөрчлөлт хийсэн өгөгдөл	33
5.5	Өгөгдлийн статистик-1	34
5.6	Өгөгдлийн статистик-2	35
5.7	Өгөгдлийн статистик-3	36
5.8	Өгөгдлийн статистик-4	37
5.9	Өгөгдлийн статистик-5	38
A.1	Бакалаврын судалгааны ажлын үечилсэн төлөвлөгөө	41
B.1	Фолдерийн бүтэц	42

ХҮСНЭГТИЙН ЖАГСААЛТ

4.1	advertisement хүснэгт	21
4.2	category хүснэгт	23
4.3	location хүснэгт	23
4.4	contactInfo хүснэгт	24

Кодын жагсаалт

2.1	Python энгийн жишээ	9
2.2	BeautifulSoup жишээ өгөгдөл цуглуулалт	10
5.1	Data Link crawling	27
5.2	Өгөгдөл цуглуулах	28
5.3	Хуудаслалтыг задлах	29
5.4	CSV файлууд хадгалах	30
5.5	Өгөгдөл цэвэрлэх функц	32
B.1	Бүх өгөгдлийг цуглуулах - dataScraping.py	42
B.2	Нэг зарын өгөгдлийг цуглуулах - adScrape.py	44
B.3	Өгөгдлийн төрөл - classTypes.py	47
B.4	Scrape хийх функц - scrape.py	47
B.5	Өгөгдөл цэвэрлэх - dataClean.py	48

УДИРТГАЛ

Мэдээллийн технологи эрчимтэй хөгжиж буй өнөөгийн нийгэмд байгууллага үйл ажиллагаа явуулж эхэлсэн цагаасаа эхлэн өгөгдлийг үйлдвэрлэсээр байдаг. Тэдгээр өгөгдлийг байнга хадгалах нь өгөгдлийн сангийн нөөцөд хортой байдаг тул өгөгдөлд шинжилгээ хийж, тэдгээрээс шаардлагатай өгөгдлүүдийг түүвэрлэн хадгалах нь чухал юм.

Өнөөдөр бид дэлхий нийтээрээ хурдтай амьдралын хэмнэлд ажиллаж, амьдарч байна. Мөн зах зээлийн хөгжил, ажил олгогчийн эрэлт хэрэгцээ ажил хайгчийн хүсэл онирхлыг оновчтой бөгөөд хурдан холбож өгөх нь нэн шаардлагатай. Өнөөгийн байдлаар энэ эрэлт хэрэгцээг хангасан тодорхой шийдвэрлэсэн мэдээллийн систем хомс байна. Иймд энэхүү бакалаврын судалгааны ажлаар ажил олгогч болон ажил идэвхтэй хайгч хоёрыг түргэн шуурхай холбож өгөх чатбот системийг хөгжүүлж байна.

1. СЭДВИЙН ТАНИЛЦУУЛГА

1.1 Оршил

Энэхүү бакалаврын судалгааны ажлын хүрээнд “Ажил олгогчдын өгөгдлийн анализ систем дээр суурилсан чатбот” сэдвийн дагуу ажил хайгчдыг ажлын байрны мэдээллээр хангах чатбот системийг хөгжүүлнэ. Ажлын байрны мэдээллийг Data Scraping аргын тусламжтайгаар, системд шаардлагатай мэдээллийг өгөгдлийн сангийн хэлбэрт оруулан бүтэцтэйгээр нэгтгэн авах бөгөөд үүнээс ажил хайгчдын дунд байдаг түгээмэл асуултуудын хариултыг өгнө. Мөн энэ системд машин сургалтын арга болох Language Understanding-ийг ашиглан хэрэглэгчийн асуултыг таамаглаж оновчтой хариулт өгөх боломжийг олгох юм.

1.2 Зорилго

Ажлын хайгчдын хэрэгцээт асуултад хариулж, ажлын байрны хүртээмжийг нэмэгдүүлэхэд энэхүү системийн гол зорилго оршино.

1.3 Зорилт

Дээрх зорилгод хүрэхийн тулд дараах зорилтуудыг тавьсан. Үүнд:

- Ашиглагдах технологиудыг сонгох, судлах
- Ижил төстэй системийн судалгаа хийх
- Системийн шинжилгээ хийх
- Системийг зохиомжлох
- Системийг хөгжүүлэх, сайжруулалт хийх

1.4 Алсын хараа

Ажлын байрны дэлгэрэнгүй мэдээллийг цуглуулснаар цаашид тэдгээрт шинжилгээ хийж хамгийн их эрэлттэй, өндөр цалинтай ажлын байр гэх зэрэг мэдээллүүдийг систем хэрэглэгчдэд хүргэх боломжтой юм.

2. СИСТЕМИЙН СУДАЛГАА

2.1 Системийн судалгаа

Сонгосон сэдэв болох “Ажил олгогчдын өгөгдлийн анализ систем дээр суурилсан чатбот” сэдвийн хүрээнд судалгаа хийхдээ чатбот системийн талаар болон өгөгдөл цуглуулгын аргын талаар судалсан. Үүний дараа ижил төстэй системийн болон ашиглагдах технологийн талаар судалгааг хийсэн болно.

2.1.1 Чатбот систем

Чатбот систем нь ихэвчлэн хэрэглэгчийн асуултыг хиймэл оюун ухааны тусламжтайгаар ойлгож, хариултыг автоматжуулах үндсэн зорилготой компьютерийн програм хангамж юм. Орчин үед хэрэглэгчдэд туслах үндсэн үүргийн дагуу чатбот системийг байгууллагууд олон янзаар ашиглах болсон. Тэдгээрээс дурдвал,

- Цэс дээр суурилсан чатбот (Menu-based chatbot)
- Түлхүүр үгийг танихад суурилсан чатбот (Keyword recognition-based chatbot)
- Машин сургалтын чатбот (Machine learning chatbot)

Цэс дээр суурилсан чатбот

Өнөөгийн зах зээлд хэрэгжиж буй чатботуудын хамгийн энгийн бөгөөд түгээмэл хэлбэр юм.[1]

¹ Хэрэглэгчийн асууж болох асуултуудыг урьдаас таамаглан хариултуудыг мод хэлбэртэйгээр бүтэцлэн хадгалдаг. Хэрэглэгч хүссэн хариултаа авахын тулд системийн хадгалсан хариултаар аялах хэрэгтэй болдог. Бусад чатботтой харьцуулбал, хариулт хязгаарлагдмал бөгөөд хэрэглэгчээс олон асуулт асууж цаг их шаарддагаараа сул талтай байдаг.

¹<https://www.engati.com/blog/types-of-chatbots-and-their-applications>

Түлхүүр үгийг танихад суурилсан чатбот

Энэхүү чатбот нь хэрэглэгчийн бичсэнийг уншиж тохиромжтой хариултыг өгдөг. Ингэхдээ өгүүлбэрийг хиймэл оюун ухааны нэг хэсэг болох эх хэлний боловсруулалт (Natural Language Processing)-ын тусламжтайгаар шинжилж түлхүүр үгийг таньж хариултыг өгдөг. Ижил төстэй олон асуултад хариулах эсвэл түлхүүр үг дутуу үед амжилтгүй болдог. Мөн хэрэглэгч хүссэн хариултаа олж чадахгүй байх болон үр дүн муутай хариулт өгсөн тохиолдолд цэс дээр суурилсан чатботыг хослуулан ашиглах нь найдвартай болдог бөгөөд түгээмэл шийдлүүдийн нэг байдаг.

Машин сургалтын чатбот

Энэ төрлийн чатбот нь өмнө хэрэглэгчийн харилцан яриан дээр хиймэл оюун ухаан болон машин сургалтын тусламжтайгаар шинжилгээ хийж, хэрэглэгчийн зан төлөв, асуултын хэв маягаас суралцдаг. Ингэснээрээ чатботод хэрэглэгчийн зарцуулах цаг эрчимтэйгээр буурах буюу хариултаа авах алхам багасгах ба хэрэглэгчийн туршлага (UX) нь түүнийгээ даган өсөх нь энэхүү чатботын үндсэн зорилго болно.

Чатботыг сонгох

Машин сургалтын чатбот нь илүү уян хатан хэрэглэгчдэд ээлтэй чатботыг бий болгодог боловч хөгжүүлэхэд цаг хугацаа их шаардагдах ба машин өөрөө суралцахад мөн хугацаа шаардагддаг. Иймд системийн нөөц, шаардлагыг харгалзан үзэж энэхүү судалгааны ажлаар түлхүүр үг танихад суурилсан чатботыг хэрэгжүүлэхийг зориод байна.

2.1.2 Өгөгдөл цуглуулгын арга

Өгөгдөл цуглуулах (data scraping) нь хэрэглэгчдэд харагдаж буй өгөгдлийг олон янзын сувгаас цуглуулан хувийн орчинд хадгалан цаашид ашиглах боломжийг олгодог хамгийн үр дүнтэй автомат өгөгдөл олборлох арга юм. Ихэвчлэн өгөгдөл цуглуулах арга нь вэбсайтаас шаардлагатай өгөгдлийг цуглуулахад ашигладаг. Өгөгдөл цуглуулж буй хүнээс хамааран олборлосон өгөгдлийг таслалаар тусгаарлагдсан утгын (Comma-Separated Values) файл эсвэл

өгөгдлийн санд хадгалах боломжтой бөгөөд нэгэнт цуглуулсан их хэмжээний өгөгдөлд судалгаа шинжилгээ хийх, худалдаа, борлуулалтын хэрэгсэл болгох зэрэг олон төрлийн боломжийг олгодог.

Вебсайтаас өгөгдлийг олборлох хамгийн түгээмэл арга нь HTML parsing буюу HTML-ийг задлан шинжлэх юм. Энэ нь вебсайтын HTML болох сайтын үндсэн бүтцийг агуулгынх нь хамтаар хуулах бөгөөд авах гэж буй өгөгдлийн зан төрхийг нь зааж өгснөөр доторх агуулгыг хамгийн хялбар бөгөөд автомат байдлаар цуглуулдаг юм. Цуглуулга хийх 2 үндсэн арга байдаг. Үүнд:

- Өгөгдлийг цуглуулж, задлах (Data scraping)
- Өгөгдлийг олж илрүүлж, хаягийг цуглуулах (Data crawling)

Өгөгдлийг цуглуулж, задлах

Нэг үгээр хэлбэл өгөгдлийг цуглуулж, задлах нь зааж өгсөн хаягийн дагуу шаардлагатай өгөгдлийг задалж, хэрэгтэй агуулгыг хөгжүүлэгчдэд өгдөг бөгөөд хүссэн өгөгдлөө задлан авах боломжийг олгодгоороо давуу талтай. Өөрөөр хэлбэл өгөгдөл олборлох програм нь зорилго буюу даалгавраа мэдэж байгаа юм.

Өгөгдлийг олж илрүүлж, хаягийг цуглуулах

Энэхүү аргачлал нь хаяг тодорхой бус үед түүнийг олж илрүүлж шаардлагын дагуу хаягийг, зарим тохиолдолд өгөгдлийг цуглуулдаг. Системийн шаардлагын дагуу өгөгдлийг цуглуулах үед хаяг алгасах, дутуу өгөгдөл цуглуулахаас сэргийлдэг давуу талтай.

Ихэвчлэн энэхүү хоёр аргыг хослуулан ашигладаг бөгөөд шаардлагад нийцэх өгөгдлийг үлдээлгүй бүгдийг нь олоход *data crawling*-ийг ашиглах бол олсон өгөгдлийг задалж, шинжлэн өгөгдлийн санд хадгалах үйлдлийг *data scraping* хийдэг. Жишээлбэл, худалдааны сайтын бараа бүтээгдэхүүний өгөгдлийг цуглуулах гэж байгаа гэж үзвэл, барааны ангиллын хаягуудыг өөрчлөгдөх бүрд хадгалан өгөгдлийг цуглуулна. Өөрөөр хэлбэл нэг нь өөрчлөлт гарахыг ажиглаж вебсайтаар мөлхөж байх бол нөгөө нь шаардлагын дагуу бүх хэрэгтэй өгөгдлийг хэдийн цуглуулсан

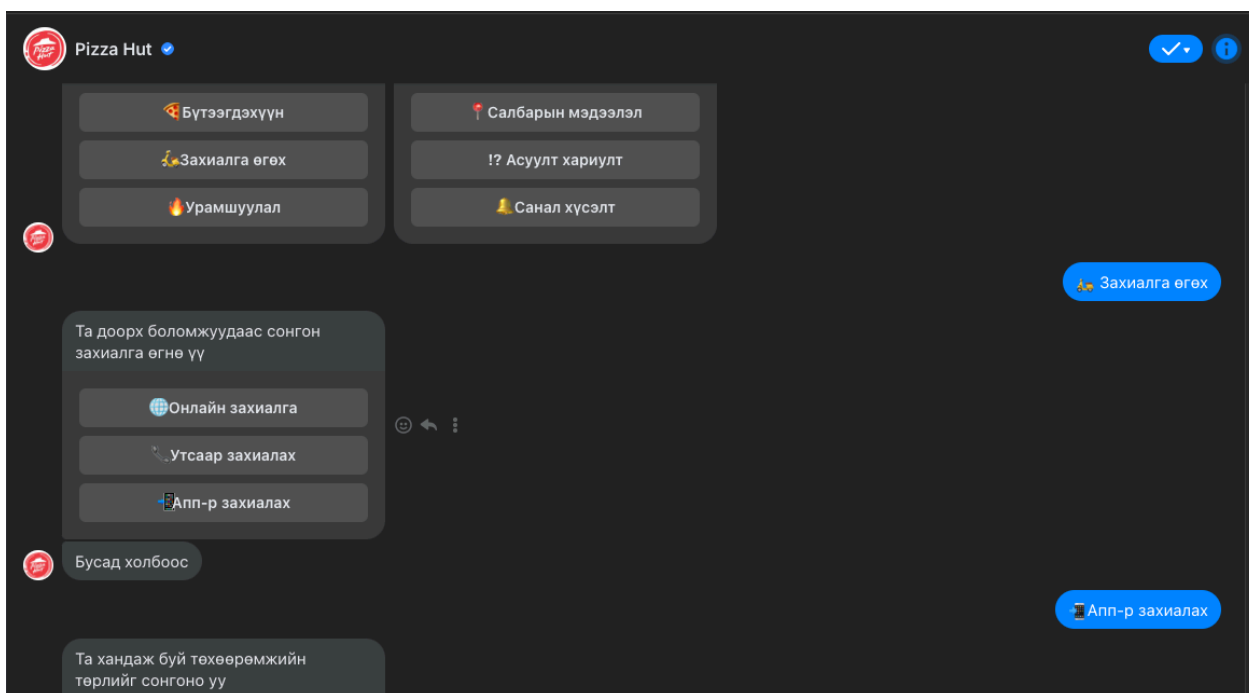
байна. Энэхүү бакалаврын судалгааны ажлын хүрээнд өгөгдлийг CSV файл үүсгэн хадгалж цаашид ашигласан болно.

2.2 Ижил төстэй системүүд

Гадаад ба дотоодын байгууллагуудын үйл ажиллагаандаа хэрэгжүүлдэг чатбот системүүдээс, *Domino's Pizza & Pizza Hut* болон *WHO's Chat bot* гэсэн гурван чатботыг сонгон авч судалгаанд оруулав.

2.2.1 *Domino's Pizza & Pizza Hut*

Domino's pizza хоолны газар нь захиалгын алхмаас эхлээд бүх мэдээллийг ганцхан *Facebook messenger chatbot* хангадаг. Чатбот эрчээ авч эхэлсэн шалтгаан нь хүмүүс, бусад хүмүүсийг хүлээлгүйгээр үйлчилгээ авах, тусламж авах зэрэг үйлчилгээг зэрэг нэвтрүүлсэнтэй холбоотой билээ. Үүний нэгэн адилаар Монголд үйл ажиллагаа явуулж буй *Pizza Hut Mongolia* юм.

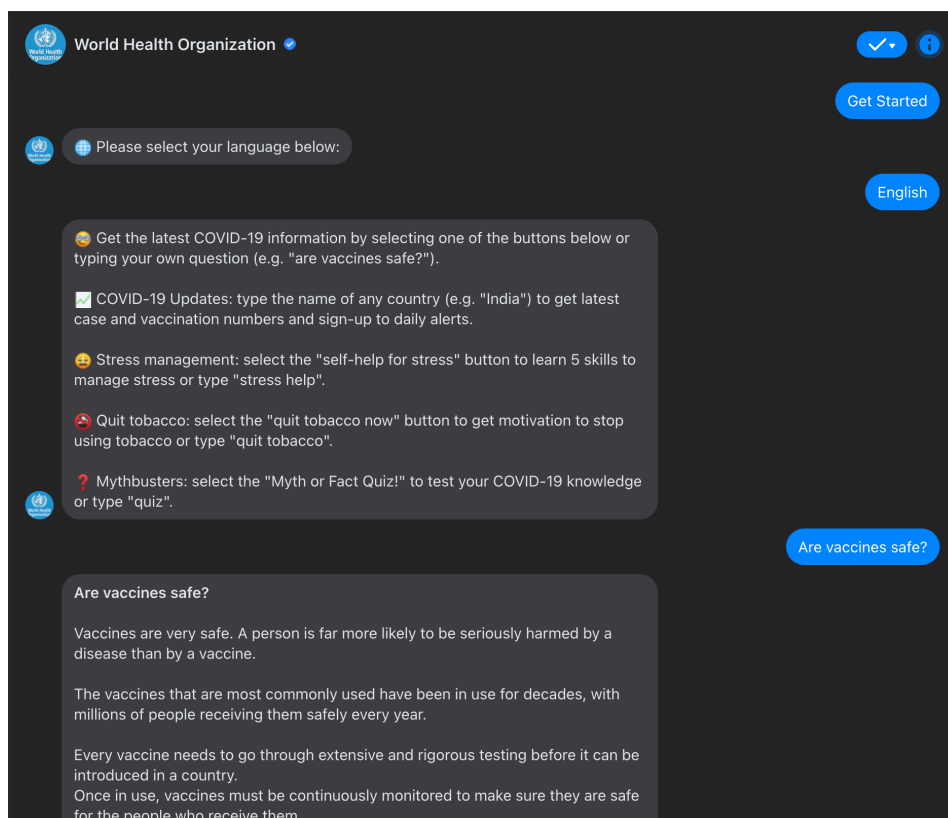


Зураг 2.1: Pizza Hut chat bot

Үйлчлүүлэгчдийн захиалга хүлээх хугацааг багасгахын тулд захиалгын үйл явцыг хурдасгаснаар тодорхой хэмжээнд нөлөөлж байгаа нь дээрх 2 жишээнээс харагдаж байна.

2.2.2 World Health Organization's Chat bot

Цар тахал болох коронавирусын эрчимтэй тархаж байх үед дэлхийн өнцөг булан бүрд оршин суугаа хүмүүст цар тахлын мэдээлэл, урьдчилан сэргийлэх арга, баталгаатай эх сурвалжийн мэдээллээр хангах зорилготой чатбот юм. Дэлхий нийтээр вакцинжуулалтын хөдөлгөөн өрнөж байх үеэр вакцины талаарх мэдээлэл, архаг хууч өвчинд нөлөөлөх талаар найдвартай, хамгийн сүүлийн үеийн албан ёсны мэдээллийг өгдөг. Хэдий халдварын тоо буурч, нийгэм өөрөө дасан зохицож байгаа хэдий ч Дэлхийн Эрүүл Мэндийн байгууллага үүргээ гүйцэтгэж чухал эх сурвалжаар хангасаар байгаагийн шинж юм.



Зураг 2.2: WHO's chat bot

2.3 Технологийн судалгаа

АЖИЛ ОЛГОГЧДЫН ӨГӨГДЛИЙН АНАЛИЗ СИСТЕМ ДЭЭР СУУРИЛСАН ЧАТ БОТ-ыг хөгжүүлэхдээ өгөгдөл цуглуулгыг *python* хэлний сан болох *BeautifulSoup* HTML өгөгдөл задлах технологийг ашигласан бөгөөд чатбот системийн түлхүүр үг таних технологийг *Python* хэлний *Framework* болох *SentenceTransformers*-ийг сонгон хөгжүүлэлтийг хийсэн. Харин цуглуулсан өгөгдлийг *CSV* файлд хадгалан, *Microsoft Bot Framework*-ийг чатбот хөгжүүлэлтэд ашиглан судалгааг дараах байдлаар хийсэн болно.

2.3.1 *Python*

Python нь дээд түвшний маш олон төрлийн програмчлалыг өөртөө шингээсэн хэл юм. Хэлний сан болон *framework*-үүд нь тасралтгүй сайжирч, шинэчлэгдэж байдаг тул бүх л төрлийн програмчлалын аргуудыг гүйцэтгэж болдог. Орчин үед машин сургалт, хиймэл оюун ухаан болон эх хэлний боловсруулалтад(NLP) түгээмэл ашигладаг болсон бөгөөд веб хүртэл хийх боломжтойгоороо давуу талтай юм. Үүнээс гадна анхлан суралцаж буй хүмүүст ойлгоход хялбар *syntax*-ийн дүрэмтэй байдаг тул хэрэглэгчдийн тоо нь javascript, java хэлүүдтэй өрсөлдөхүйц байдаг.



Зураг 2.3: Python лого

```
1 x = 5
2 name = 'Sainzolboo'
3 print(x)
4 print(name)
```

Код 2.1: Python энгийн жишээ

Python програмчлалын хэл нь ойлгоход маш хялбар бөгөөд өөр дээрх функцууд нь шууд утгаараа ойлгомжтой байдаг. Syntax-ийн хувьд ; ашигладаггүй ба догол мөрөөр програмчлалын үндсэн схемийг гаргадгаараа онцлог хэл юм.

2.3.2 *BeautifulSoup*

Өгөгдөл цуглуулгын олон технологиудын нэг нь *BeautifulSoup* бөгөөд *python* програмчлалын хэлний сан юм. Энэ нь өгсөн вебсайтын хаяг (Url)-ийн дагуу бүх HTML өгөгдлийг агуулгын хамтаар нь хэрэглэгчид өгдөг. HTML хэл нь мод хэлбэртэй байдаг бөгөөд түүний хүүхэд элементүүдийн агуулгыг шаардлага болон түлхүүр үгийн дагуу цуглуулах зарчмаар ажилладаг.



Зураг 2.4: BeautifulSoup лого

Бакалаврын судалгааны ажлын сэдвийн дагуу ажлын байр олгогчдын мэдээлэл болон ажлын байрны мэдээллийг **zangia.mn**-ээс *BeautifulSoup* ашиглан цуглуулсан. Доорх кодын жишээнд бүх ажлын байрны ангилал болон шүүлтүүрийн агуулгыг цуглуулсан бөгөөд жишээнд зориулж зөвхөн эхний ангиллын мэдээллийг харуулав.

```
1 from bs4 import BeautifulSoup
2 import requests
3 from urllib.error import HTTPError
4
5 url = 'https://zangia.mn/'
```

```

6  try:
7      response = requests.get(url)
8      response.raise_for_status()
9  except HTTPError as error:
10     print(error)
11  soup = BeautifulSoup(response.text, 'html.parser')
12  navigatorList = soup.find_all('div', class_='filter')
13  print(navigatorList[0])

```

Код 2.2: BeautifulSoup жишээ өгөгдөл цуглуулалт

```

<div class="filter">
  <h3>
    Онцлох
  </h3>
  <div>
    <a href="job/list/x.1">
      Удирдах албан тушаалын ажлын байр
    </a>
  </div>
  <div>
    <a href="job/list/x.3">
      Англи хэлний 100%-н мэдлэг шаардах ажлын байр
    </a>
  </div>
  <div>
    <a href="job/list/x.2">
      Ажлын туршлага шаардахгүй ажлын байр
    </a>
  </div>

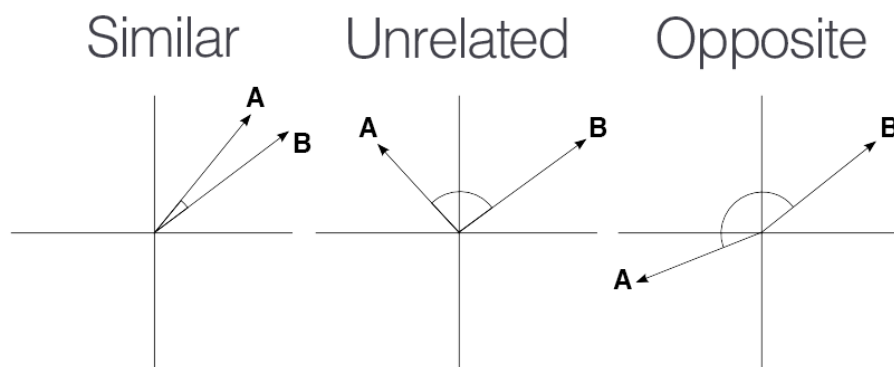
```

Зураг 2.5: Өгөгдөл цуглуулалтын жишээний үр дүн

2.3.3 SentenceTransformers

Python хэлний framework болох SentenceTransformers [2] буюу өгүүлбэр хувиргалт нь өгүүлбэр болон текстийн ижил төстэй байдал болон утгын хувьд адил байдлыг *cosine-similarity*² -ийн тусламжтайгаар тооцоолдог. Энэхүү тооцооллыг цаашид өгүүлбэрийн ижил төстэй байдлыг харьцуулах, хайлт хийх, түүнд шинжилгээ хийх зэргээр ашиглаж болно. Доорх зурагт өгүүлбэрт хувиргалт хийж, шинжилгээний үр дүнгийн вектор хоорондын өнцгөөр хэрхэн тодорхойлогддог болох талаар харуулав.

²Cosine-similarity нь өгөгдлийн шинжилгээнд 2 тооны ижил төстэй байдлыг вектор үржвэрээр илэрхийлдэг.



Зураг 2.6: Cosine similarity утгын график

SentenceTransformers-ийг дэлхийн 100 гаруй хэл дээр урьдчилан бэлтгэн, сургасан эх хэлний боловсруулалт (NLP)-ын загваруудыг ашиглаж болдгоороо давуу талтай.



Зураг 2.7: SBERT.net лого

Чатбот системийн хувьд монгол хэлийг танин ашиглах боломжтой загвар болох *distiluse-base-multilingual-cased-v2[3]*-ийг ашиглан хийж гүйцэтгэв.

Хоёр өгүүлбэрийг *cosine-similarity* ашиглан ижил төстэй байдлыг илэрхийлэх жишээг доор харууллаа. Эх кодыг utf-8 формат танихгүй байсан тул зураг хэлбэрээр орууллаа.

```
tests.py > ...
1  from sentence_transformers import SentenceTransformer, util
2  model = SentenceTransformer(
3      'sentence-transformers/distiluse-base-multilingual-cased-v2')
4
5  emb1 = model.encode("Энэ бол улаан малгайтай муур юм.")
6  emb2 = model.encode("Энэ бол миний улаан малгайтай нохой.")
7
8  cos_sim = util.cos_sim(emb1, emb2)
9  print("Cosine-Similarity:", cos_sim)
10

PROBLEMS 14 OUTPUT DEBUG CONSOLE TERMINAL JUPYTER

/usr/local/bin/python3 /Users/zolboo/Desktop/bachelor/employmentAnalysis/project/sbert.py
zolboo@Sainzolboos-MacBook-Pro project % /usr/local/bin/python3 /Users/zolboo/Desktop/bachelor/employmentAna
Cosine-Similarity: tensor([[0.7287]])
zolboo@Sainzolboos-MacBook-Pro project %
```

Зураг 2.8: cosine-similarity ашигласан жишээ

2.3.4 Linear Regression - Шугаман Регресс

Linear Regression буюу шугамар регрессийн загвар нь шулуун шугамыг ашигладаг бол логик болон шугаман бус регрессийн загвар нь муруй шугамыг ашигладаг. Шугаман регресс нь бие даасан хувьсагч хэрхэн өөрчлөгдөхийг тооцоолох боломжийг олгодог. Хоёр тоон хувьсагчийн хоорондын хамаарлыг тооцоолоход шугаман регрессийг ашигладаг бөгөөд ихэвчлэн дараах нөхцөлд ашиглагддаг. Үүнд:

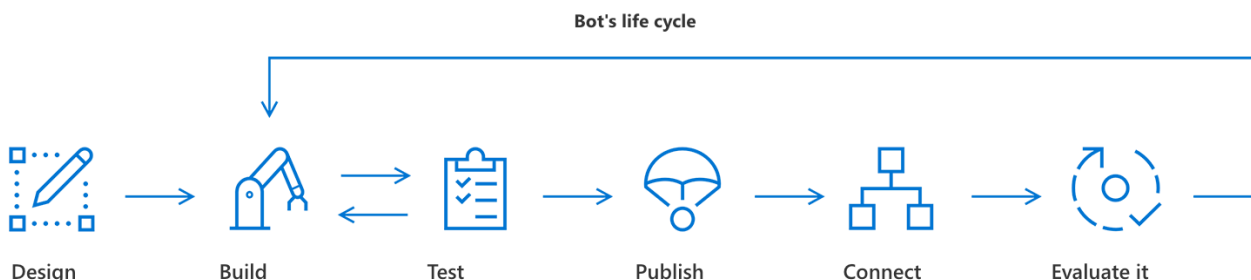
- Хоёр хувьсагчийн хоорондын нөлөөлөл
- Бие даасан хувьсагчийн тодорхой утга дахь хамааралтай хувьсагчийн утга
- Вариацийн нэгэн төрлийн байдал
- Ажиглалтын бие даасан байдал
- Нормал байдал

2.3.5 Comma Separated Values - CSV файл

CSV нь өгөгдлийн утгуудыг тусгаарлахад таслал ашигладаг текст файл юм. Файлын мөр бүр нь өгөгдөл байдаг бөгөөд харгалзах утгуудад текст файлыг бичих энгийн өгөгдөл хадгалах технологи юм. Их өгөгдөлтэй хялбар харьцах боломжийг олгодгоороо давуу талтай.

2.3.6 Microsoft Bot Framework

*Bot Framework*³ нь Microsoft-ийн *Azure Bot Service*-ийн тусламжтайгаар чатботыг турших, үүсгэх, удирдах, хэрэглээнд нэвтрүүлэх гэх мэт боломжуудыг нэг дор хангаж өгдөг. Энэхүү боломжуудын хүрээнд асуулт хариултыг зохицуулах, хэрэглэгчид зориулсан User Interface бүтээх, Language Understanding аргыг ашиглах гэх мэт үйлдлүүдийг хийх боломжтой. Bot бүтээх үйл явцыг Azure Bot Service болон Bot Framework нь ихэд хөнгөвчилж өгдөг бөгөөд доорх зурагт үзүүлсэн дарааллын дагуу Bot системийг бүтээдэг.



Зураг 2.9: Чатботын амьралын мөчлөг

Design

Design буюу загварчлах нь төслийн төлөвлөгөөг гаргах юм. Өөрөөр хэлбэл, системийн зорилго, үйл явц, хэрэглэгчийн хэрэгцээг сайтар судлах нь амжилттай Bot систем бүтээх чухал хэсэг юм.

Build

Бот системийг угсрах буюу хөгжүүлэх үйл явц юм. Энэ алхамд хөгжүүлэгч хэрэглэгчийн

³<https://dev.botframework.com/>

харагдах хэсгийг загварчлах бөгөөд хөгжүүлэлтийн орчин нь *Azure Portal*, JavaScript, Python болон C програмчлалын хэлнүүдээс сонгож хөгжүүлэлтийг гүйцэтгэх явц юм. Мөн системийн шаардлагыг тодорхойлсны дагуу бот системийг өргөжүүлж ашиглах боломжтой бөгөөд тэдгээрээс дурдвал:

- Эх хэлний боловсруулалт (NLP)
- Асуулт хариултыг сайжруулан мэдлэгийн сан үүсгэх
- Хэрэглэгчийн интерфэйсийг сайжруулах

Test

Програм хангамжийн хөгжүүлэлтийн амьдралын мөчлөгийн адилаар тестийн үйл явцыг алгасаж болохгүй. Нэгэнт хэрэглэгчийн гарт бот системийг оруулахаас өмнө гарч болох алдаа дутагдлыг засан сайжруулах шаардлагатай. Иймд Bot системийг publish хийхээс өмнө заавал туршиж үзэх шаардлагатай. Энд Microsoft-ийн өөрсдийнх нь бие даасан програм болох *Bot emulator*-ийг ашиглан хөгжүүлэлтийн орчинд туршиж үзэх боломжийг олгодог.

Publish

Тестийн шатны дараа бот систем ашиглахад бэлэн болсон гэж үзсэн үед төсөл эсвэл чатботыг олон нийтэд ил болгох явдал юм.

Connect

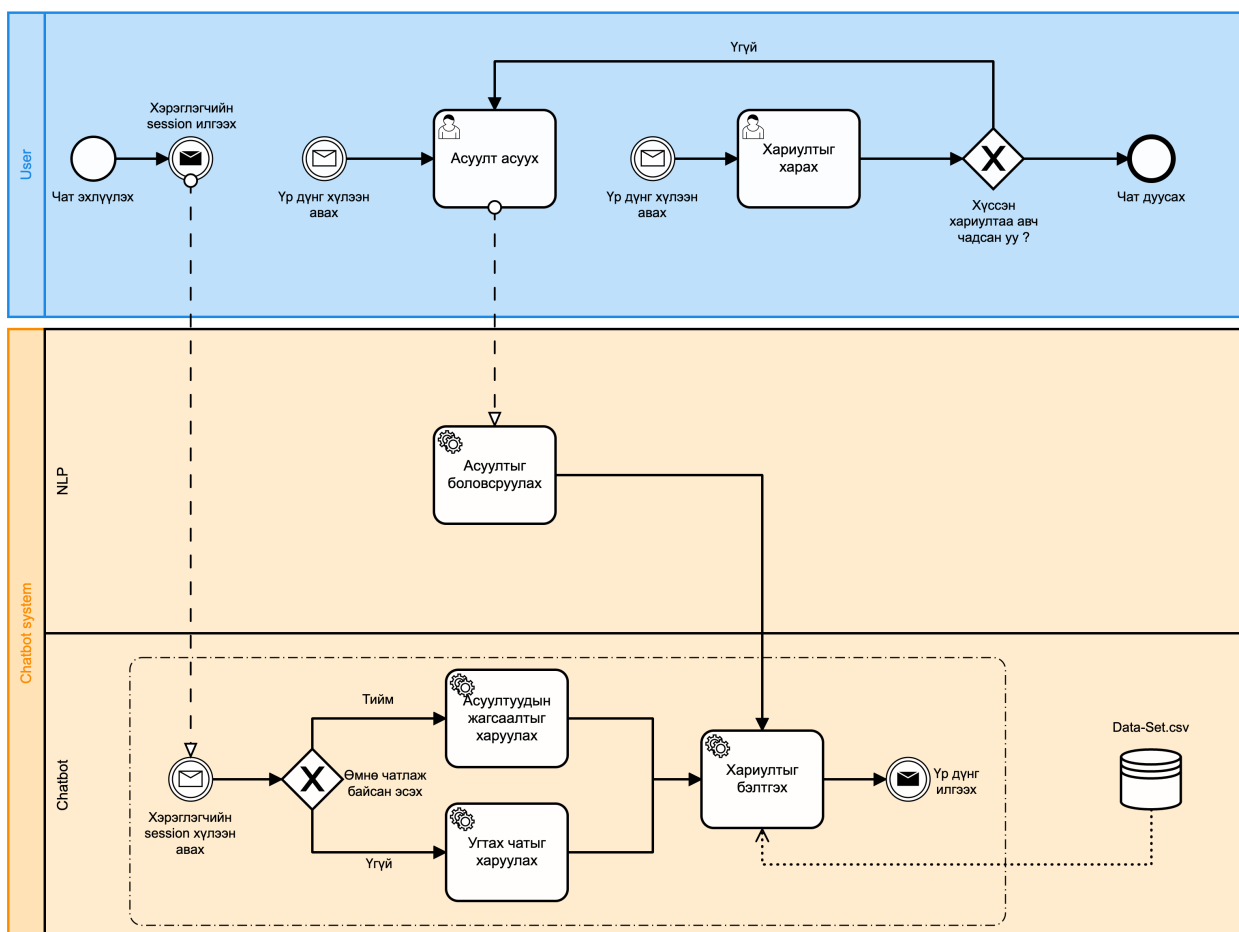
Bot системээ Facebook messenger, Microsoft Teams, Telegram, Skype гэх мэт чат сувгуудыг өөрийн Bot-той холбоно.

Ингэж бүх мөчлөгийг дууссаны дараа хөгжүүлэгч хэрэглэгчийн ашиглаж буй байдал дээр анализ хийж системийг дахин сайжруулах боломжтой бөгөөд буцаад угсрах үйл явцруу шилжих юм.

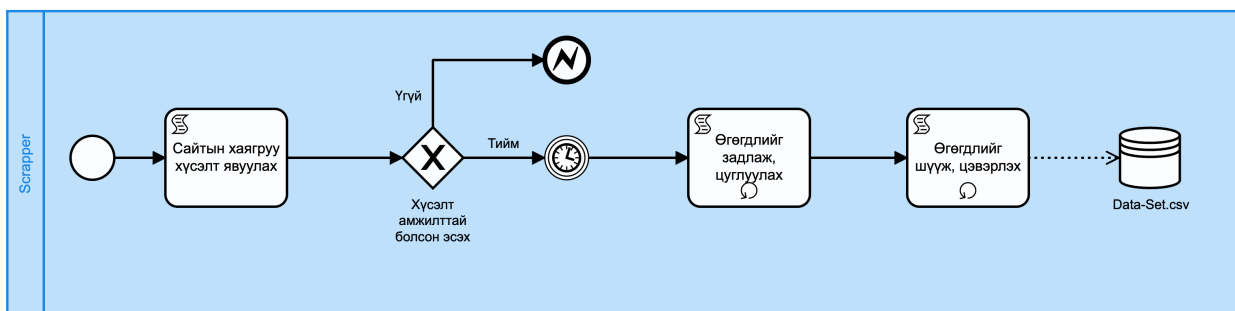
3. СИСТЕМИЙН ШИНЖИЛГЭЭ

3.1 Бизнесийн үйл ажиллагааны шинжилгээ

Бизнес процессийн модель нь чатбот системийн үндсэн процесс буюу үйл ажиллагааны явцыг BPMN-2.0 ашиглан дүрслэн харуулав [5]. Диаграммд дүрслэхдээ оролцогч талууд болох системүүдийг тус тусын *pool* дотор дүрсэлсэн бол дэд процесс буюу *subprocess*-ийг *lane*-д дүрсэлж хоорондын хамаарлыг харууллаа.



Зураг 3.1: BPMN-1



Зураг 3.2: BPMN-2

3.2 Хэрэглэгч

Чатбот системийг ямар ч хүн хэрэглэх боломжтой бөгөөд олон нийтэд нээлттэй байна. Системийн гол зорилго нь ажил хайж буй хэрэглэгчдэд ажлын байрны цогц мэдээллийг олгох зорилготой байх тул хэрэглэгчдийг дараах байдлаар тодорхойлж болно. Үүнд:

- Ажлын байр хайж буй хүн
- Хөгжүүлэгч

3.3 Функционал шаардлага

Дараах хэсэгт чатбот системд тавигдах функционал шаардлагуудыг харуулсан болно.

ФШ 1 Чатбот нь харилцан яриа эхэлмэгц хариу өгдөг байна.

ФШ 2 Чатбот нь ямар ч оролтод хариу өгнө.

ФШ 3 Хэрэв чатбот нь оролтод хариу өгч чадхааргүй байвал бусад асуултуудыг санал болгож ойлгомжгүй утга оруулсныг илэрхийлнэ.

ФШ 4 Чатботын санал болгох асуултууд нь цэс хэлбэртэй харагдана.

ФШ 5 Чатботын цэсэн дээр нэг товшилтоор асуултын хариултыг харуулдаг байна.

ФШ 6 Алхам бүрд үндсэн цэсрүү буцах сонголтыг харуулдаг байна.

ФШ 7 Чатботны хариулт нь текстэн хэлбэрээр хэрэглэгчид харагдана.

ФШ 8 Чатбот нь зөвхөн Монголоор бичсэн асуултад хариулт өгнө.

ФШ 9 Чатбот нь дэлгэрэнгүй мэдээллийг цэс хэлбэрээр сонгуулан харуулж чаддаг байна.

3.4 Функционал бус шаардлага

Бэлэн болон найдвартай байдал (Availability & Reliability)

ФБШ 01 Чатбот систем өдрийн аль ч цагт 99.999% ажиллагаатай байх ёстой.

ФБШ 02 Ямар ч хүсэлт ирсэн чатбот 100% хариу өгдөг байна.

Гүйцэтгэлтэй байдал (Performance)

ФБШ 03 Чатботын байршуулсан сувагт, шаардлагаас хамаарч ямар ч төхөөрөмжөөс хандаж болно.

ФБШ 04 Зарим тохиолдолд чатботын гүйцэтгэл нь хэрэглэгчийн интернет болон төхөөрөмжийн үйлдлийн системийн хувилбараас хамаарч болно.

Дэмжих чадвар (Supportability)

ФБШ 05 Чатботын эх кодыг *github* дээр нээлттэй эхийн систем хэлбэрээр байршуулна.

Хэрэгцээт байдал (Usability)

ФБШ 06 Чатбот нь хэрэглэхэд хялбар, ойлгомжтой байна.

ФБШ 07 Чатботны цэс нь ойлгомжтой цөөн үгээр илэрхийлэгдсэн байна.

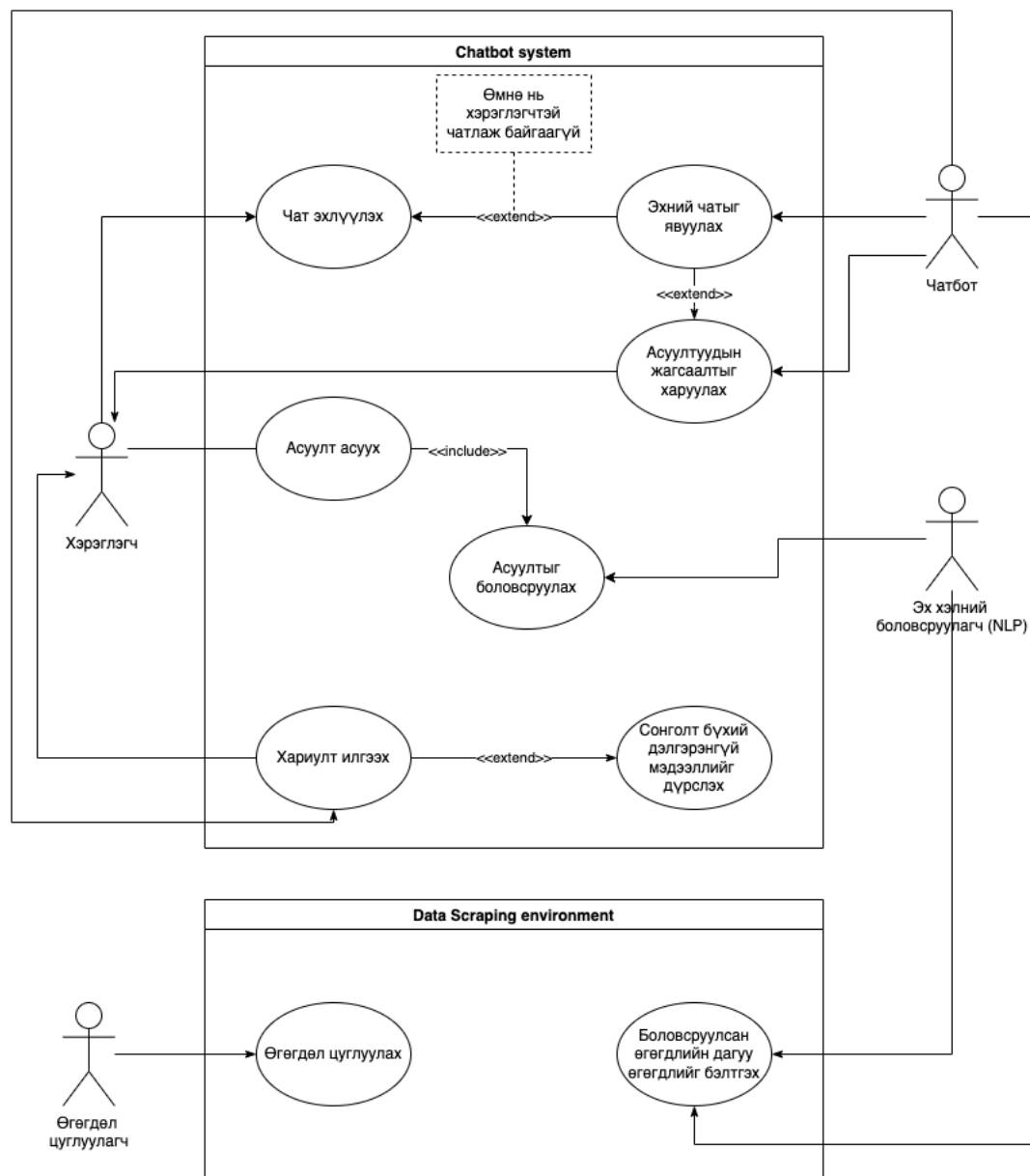
ФБШ 08 Чатботны цэсийн хэмжээ дарагдахуйц том байна.

Аюулгүй байдал (Security)

ФБШ 09 Чатбот системийн байршуулсан сувгийн стандартын дагуу хэрэглэгчийн мэдээллийг өгөгдлийн санд хадгалахгүй байна.

3.5 Use case диаграмм

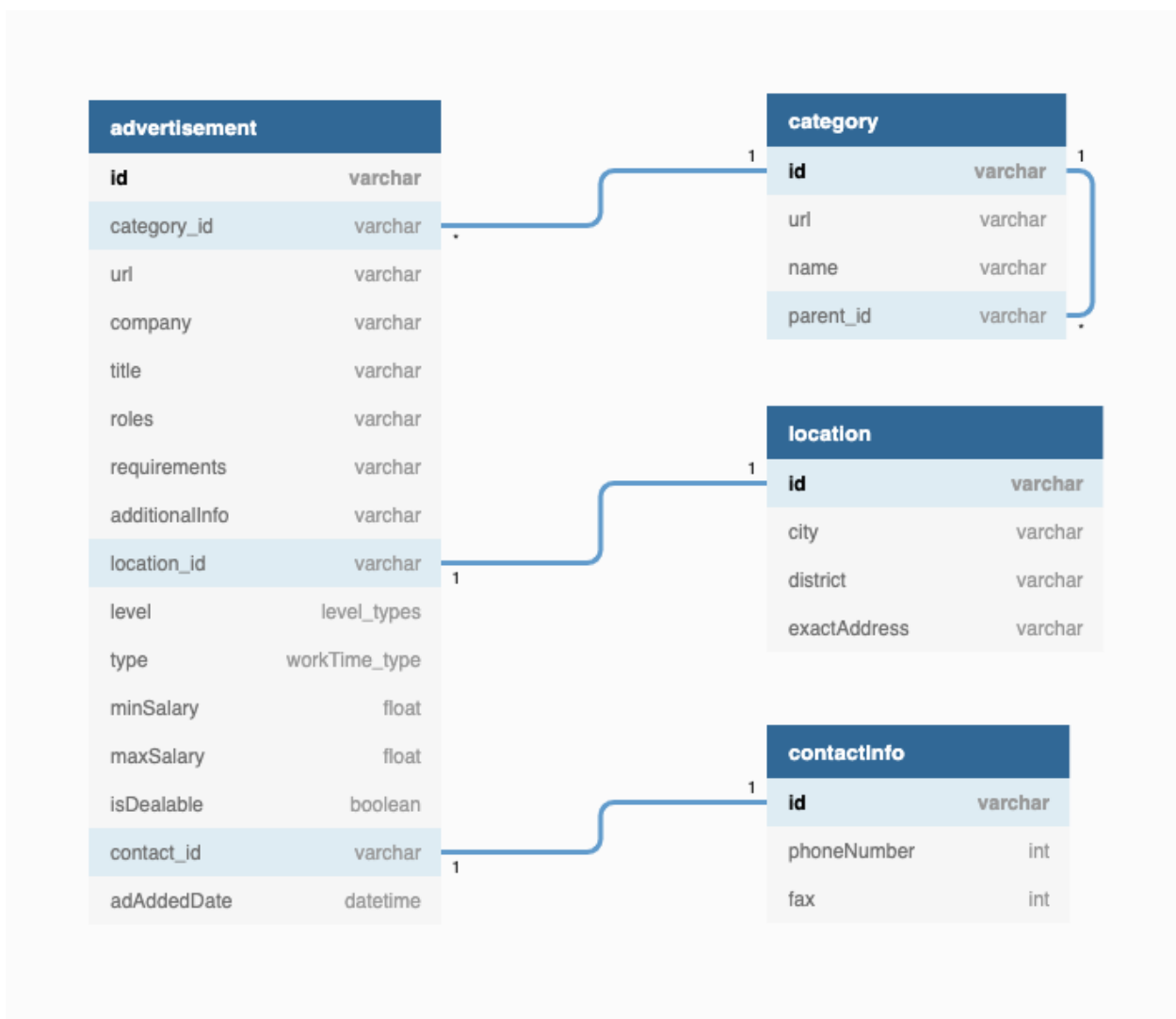
Чатбот системийн use-case диаграммыг байдлаар тодорхойлов [4].



Зураг 3.3: Use Case диаграмм

4. СИСТЕМИЙН ЗОХИОМЖ

4.1 Өгөгдлийн сангийн диаграмм



Зураг 4.1: Өгөгдлийн сангийн диаграмм

4.2 Өгөгдлийн элемент

Чатбот системийн өгөгдлийн сангийн диаграммд харуулсан хүснэгтүүдэд агуулагдах мэдээлэл болон үүргийн талаар дэлгэрэнгүй тайлбарласан болно.

4.2.1 *advertisement* - Ажлын байрны зар

Ажлын байрны зар нь ямар категори буюу ангилалд, ямар холбоо барих хаягийн хамтаар хадгалагдаж буй мэдээлэл болон бусад дэлгэрэнгүй мэдээллийг харуулсан байна.

Table 4.1: advertisement хүснэгт

№	Баганын нэр	Түлхүүр өгөгдөл	Өгөгдлийн төрөл	Хоосон утга	Тайлбар
1	id	PK	varchar	not null	Ажлын байрны зарын дахин давтагдашгүй дугаар
2	category_id	FK	varchar	not null	Ажлын байрны зард хамаарах ангиллын дугаар
3	url		varchar	not null	Ажлын байрны зарын хаяг
4	company		varchar	not null	Ажил олгогч компани / хүн
5	title		varchar	not null	Ажлын зарын гарчиг
6	roles		varchar	null	Гүйцэтгэхүндсэн үүрэг
7	requirements		varchar	null	Ажлын байранд тавигдах шаардлага
8	additionalInfo		varchar	null	Нэмэлт мэдээлэл
9	location_id	FK	varchar	not null	Ажлын байрны зард хамаарах ангиллын дугаар
10	level		level_types	null	Ажлын түвшин
11	type		workTime_type	null	Ажиллах цагийн төрөл

№	Баганын нэр	Түлхүүр өгөгдөл	Өгөгдлийн төрөл	Хоосон утга	Тайлбар
12	minSalary		float	null	Доод цалин
13	maxSalary		float	null	Дээд цалин
14	isDeable		boolean	null	Тохиролцох эсэх
15	contact_id	FK	varchar	not null	Ажлын байрны зард хамаарах холбоо барих хаягийн дугаар
16	adAddedDate		datetime	not null	Зар нийтэлсэн огноо

Энд *level* буюу ажлын түвшин, *type* буюу ажлын цагийн өгөгдлийн төрлийг тодорхойлохдоо дараах байдлаар зааж өгсөн.

Enum level_types буюу ажлын түвшний шаардлага нь дараах үндсэн 4 өгөгдлийн төрлөөс хамаарна:

- student - Оюутан / дадлагажигч
- professional - Мэргэжлийн
- occupasionDoesntRequire - Мэргэжил шаардахгүй
- intermediateManagemet - Дунд шатны удирдлага
- topLevelManagemet - Дээд шатны удирдлага

workTime_type буюу ажиллах цагийн нөхцөл нь дараах үндсэн 4 өгөгдлийн төрлөөс хамаарна:

- shift - Ээлжийн
- fullTime - Бүтэн цагийн
- halfTime - Хагас цагийн
- contract - Гэрээт / зөвлөх

- seasonal - Улирлаар

4.2.2 category - Ангилал

Ажлын байрны зарын бүх ангиллуудын хаяг болон нэрийн мэдээллийг хадгалах хүснэгт юм. Ангиллууд нь дэд ангилал байж болох учир түүнийг эцэг ангиллын дугаарыг хадгалах байдлаар зохиомжлов.

Table 4.2: category хүснэгт

№	Баганын нэр	Түлхүүр өгөгдөл	Өгөгдлийн төрөл	Хоосон Утга	Тайлбар
1	id	PK	varchar	not null	Ажлын байрны зарын ангиллын дугаар
2	url		varchar	not null	Ангиллын хаяг
3	name		varchar	not null	Ангиллын нэр
4	parent_id	FK	varchar	null	Эцэг ангиллын дугаар

4.2.3 location - Байршил

Ажлын байрны байршил болон хот, аймаг, дүүргийн дэлгэрэнгүй өгөгдлийг хадгална.

Table 4.3: location хүснэгт

№	Баганын нэр	Түлхүүр өгөгдөл	Өгөгдлийн төрөл	Хоосон Утга	Тайлбар
1	id	PK	varchar	not null	Ажлын байрны зарын хаягийн дугаар
2	city		varchar	null	Ажлын байрны зарын байрших хот, аймгийн нэр

№	Баганын нэр	Түлхүүр өгөгдөл	Өгөгдлийн төрөл	Хоосон Утга	Тайлбар
3	district		varchar	null	Ажлын байрны зарын байрших дүүрэг, сумын нэр
4	exactAddress		varchar	null	Дэлгэрэнгүй хаяг

4.2.4 contactInfo - Холбоо барих

Ажлын байр олгогчийн хаягийн дэлгэрэнгүй өгөгдлийг хадгалана.

Table 4.4: contactInfo хүснэгт

№	Баганын нэр	Түлхүүр өгөгдөл	Өгөгдлийн төрөл	Хоосон Утга	Тайлбар
1	id	PK	varchar	not null	Ажил олгогчтой холбоо барих хаягийн дугаар
2	phoneNumber		int	null	Ажил олгогчийн утасны дугаар
3	fax		int	null	Ажил олгогчийн факс дугаар

4.3 Өгөгдлийн сангийн холбоосын тайлбар

- Нэг ангилал буюу категорид олон ажлын байрны зар байж болно.
- Нэг ангилал буюу категорид олон категори байж болно.
- Нэг ажлын байрны зард нэг байршлын мэдээлэл байна.
- Нэг ажлын байрны зард нэг холбоо барих хаягийн мэдээлэл байна.

5. ХЭРЭГЖҮҮЛЭЛТ, ҮР ДҮН

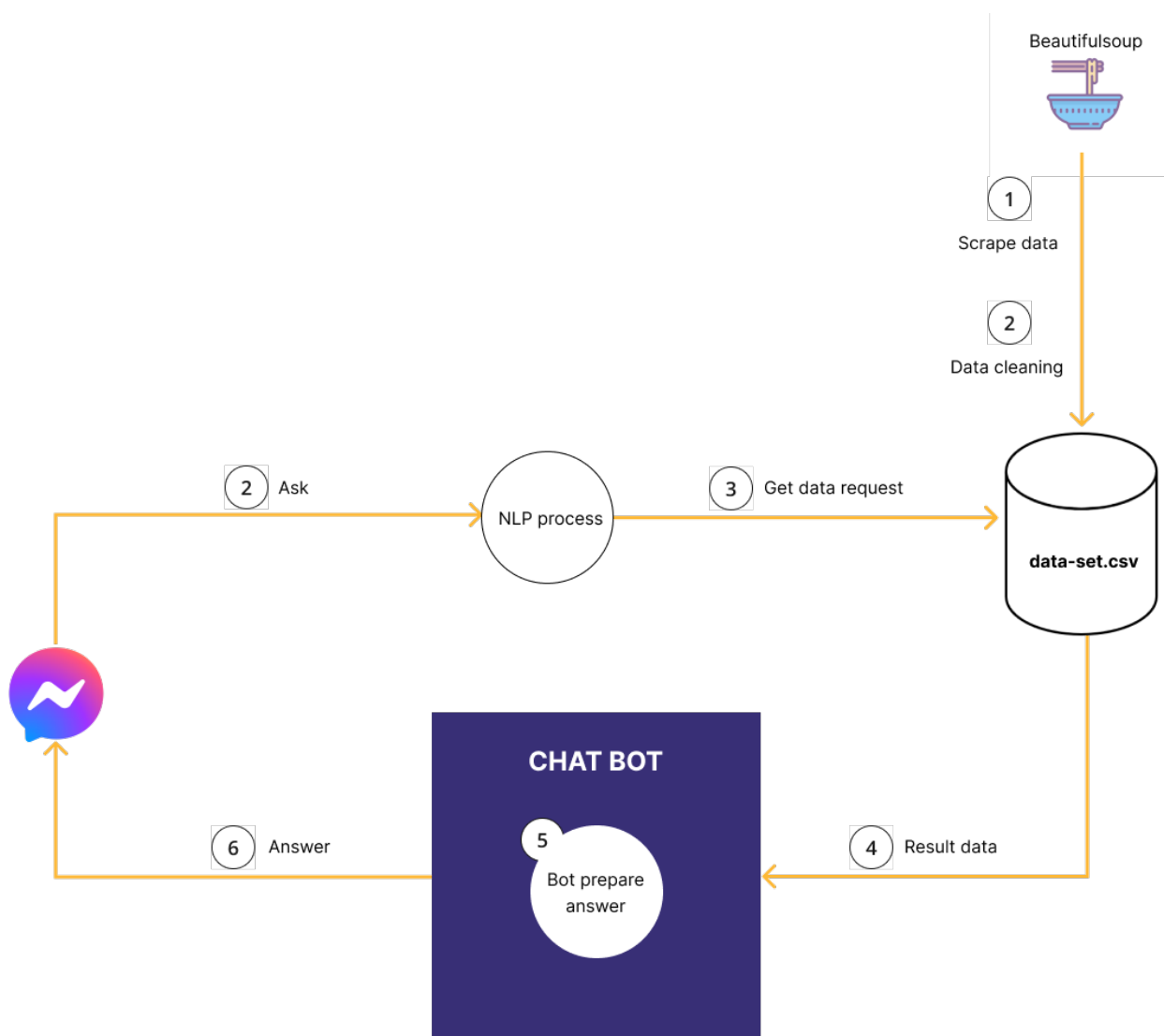
5.1 Хөгжүүлсэн байдал

Чатбот системийн хөгжүүлэлтийг хийхдээ шаардлагууд дээр үндэслэн, үечилсэн төлөвлөгөө болон шаардлагатай хөгжүүлэлтийг дэс дараалалтайгаар хийж гүйцэтгэсэн.

- Өгөгдөл цуглуулах
- Өгөгдлийг нэгтгэх, цэвэрлэх
- Системийн шаардлага, үйл ажиллагааг тодорхойлох
- Өгөгдөлд анализ хийх
- Эх хэлний боловсруулалт хийх
- Чатбот хөгжүүлэх

гэсэн дарааллын дагуу хөгжүүлэлтийг хийсэн болно.

Доорх зурагт чатбот системийн үндсэн процессийн зураглал харагдаж байна.



Зураг 5.1: Үндсэн процесс зураглал

5.1.1 Өгөгдөл цуглуулах

Үндсэн ашиглагдах өгөгдөл болох ажил олгогчид, ажлын байрны өгөгдлийг **zangia.mn**-ээс BeautifulSoup ашиглан авсан. Эхлээд вебсайтынхаа HTML бүтцийг нь судалж, авах өгөгдлийнхөө класс утгуудыг (className) олж авах нь зөв юм. Вебсайтаас өгөгдөл цуглуулах 2 үндсэн арга байдгаас өгөгдлийг олж илрүүлж, хаягийг цуглуулах (data crawling) аргаар бүх ангиллуудын хаяг (url)-уудын түүж авна. Харин data scraping нь тэр хооронд олсон бүх хаягуудаараа явж

хэрэгтэй агуулгыг цуглуулна.¹

```
1  initialUrl = 'https://www.zangia.mn/'
2  today = str(date.today())
3  # all categories set
4  categorySet = set()
5  # all advertisement's link set
6  adUrlDict = {}
7  # all ads object set
8  adsSet = set()
9
10 # scrape initial links
11 soup = useScrape(initialUrl)
12 navigatorList = soup.find_all('div', class_='filter')
13 for navigator in navigatorList:
14     if navigator.find('h3').text.strip() != 'Salbar, mergejil':
15         continue
16     # ALL CATEGORY LINKS
17     categoryList = navigator.find_all('div')
18
19 for categoryItem in categoryList:
20     categories = categoryItem.find('a')
21     url = initialUrl + categories['href']
22     tempCategory = Category(url, categories.text, '')
23     soup = useScrape(url)
24     subCategory = soup.find('div', class_='pros')
25     # ALL SUBCATEGORY LINKS
26     subCategoryList = subCategory.find_all('a')
```

¹Кодын жишээг оруулахад utf-8 формат танихгүй байсан тул монголоос галиглаж бичсэн болно.

```

27     for subCategoryItem in subCategoryList:
28         subCategoryUrl = initialUrl + subCategoryItem['href']
29         tempSubCategory = Category(
30             subCategoryUrl, subCategoryItem.text, tempCategory.name)
31         categorySet.add(tempSubCategory)

```

Код 5.1: Data Link crawling

Дээрх код нь эхлээд вебсайтруу орж “filter” класс доторх “Салбар, мэргэжил” гэсэн хэсгээс бүх эцэг категориудыг data crawling хийж авч байна. Үүний дараа хүүхэд категориудыг олж categorySet дотор бүх хаягуудыг хийж хадгалж байна.² Энд categorySet set-ийн элемент нь category төрлийн объект бөгөөд өгөгдлийн сангийн диаграм дээр тодорхойлж өгсөн байгаа. Ингэснээр data crawling-ийг зогсоож, цуглуулсан хаягаасаа өгөгдлөө цуглуулъя.

```

1  for categoryItem in categorySet:
2      if categoryItem.parentId == '':
3          continue
4      soup = useScrape(categoryItem.url)
5      hasPagination = soup.find('div', class_='page-link')
6      pagesUrl = []
7      if hasPagination != None:
8          pagesUrl = createLinkList(hasPagination, categoryItem.url)
9      else:
10         pagesUrl.append(categoryItem.url)
11     for pageUrl in pagesUrl:
12         soup = useScrape(pageUrl)
13         ads = soup.find_all('div', class_='ad')
14         # CREATE UNIQUE AD DICTIONARY
15         for ad in ads:

```

²Python хэлний set өгөгдлийн төрөл нь давхацахгүй утгуудын хүснэгт гэж хэлж болно.

```

16         adUrl = initialUrl+ad.find('a', class_=None)['href']
17         adUrlDict[adUrl] = categoryItem
18     pagesUrl.clear()

```

Код 5.2: Өгөгдөл цуглуулах

Дээрх кодоод бүх хүүхэд категориудын дотор агуулагдаж буй зарын мэдээллийг цуглуулж байна. Ингэхдээ эхлээд категори доторх өгөгдлүүд нь хуудаслагдсан (pagination) байх боломжтой бөгөөд хэрэв олон хуудастай байвал хаягуудыг нь угсарч тэдгээрээс ч мөн өгөгдлийг нь цуглуулах ёстой юм.

```

1  from .scrape import UseBeautifulSoup as useScrape
2
3  def createLinkList(pagination, url) -> array:
4      linkList = []
5      total = int(useRegex(pagination.find_all('a')[-1]['href']))
6
7      for i in range(total + 1):
8          if i == 0:
9              continue
10         link = url + '/pg.' + str(i)
11         linkList.append(link)
12     return linkList

```

Код 5.3: Хуудаслалтыг задлах

Энэ хэсэгт хуудаслан дугаарласан хэсгийн хамгийн сүүлийн тоог авч *createLinkList* функцруу дамжуулснаар тухайн категорийн бүх өгөгдлийг цуглуулах боломж үүсч байгаа юм. Ингээд дахин data crawling хийж бүх хаягуудыг цуглуулж энэ удаад dictionary үүсгэж зарын хаягуудыг хадгалсан. Энд dictionary үүсгэхдээ хаягийг нь түлхүүр (key) болгож категори объектыг нь утга(value) болгож хадгалсан. Мэдээж хэрэг dictionary нь түлхүүр давхцахаас сэргийлдэг тул

бид ямар нэгэн байдлаар нэг зарын өгөгдлийг 2 удаа цуглуулах эрсдэлгүй болж байна.³

Харин одоо үүсгэсэн dictionary-оо ашиглан өгөгдлөө CSV файлуугаа бичихэд ашиглаж болно.

```
1 for adUrl in adUrlDict:
2     print(adUrl)
3     try:
4         tempAdItem = useAdScrape(adUrl)
5         tempAdItem.setCategory(adUrlDict[adUrl])
6         file.write(
7             tempAdItem.category.parentId+'\t' +
8             tempAdItem.category.name+'\t' +
9             tempAdItem.url+'\t' +
10            tempAdItem.company+'\t' +
11            tempAdItem.title+'\t' +
12            tempAdItem.roles+'\t' +
13            tempAdItem.requirements+'\t' +
14            tempAdItem.additionalInfo+'\t' +
15            tempAdItem.city+'\t' +
16            tempAdItem.district+'\t' +
17            tempAdItem.level+'\t' +
18            tempAdItem.type+'\t' +
19            tempAdItem.minSalary+'\t' +
20            tempAdItem.maxSalary+'\t' +
21            tempAdItem.isDealable+'\t' +
22            tempAdItem.address+'\t' +
23            tempAdItem.phoneNumber+'\t' +
24            tempAdItem.fax+'\t' +
```

³Нэг зарын өгөгдлийг цуглуулахад интернетийн хурдаас хамааран 0.2-оос 0.5 секунд хугацаа зарцуулдаг

```

25         tempAdItem.adAddedDate+'\n')
26     del tempAdItem
27 except:
28     print('Ad writing error')
29 file.close()

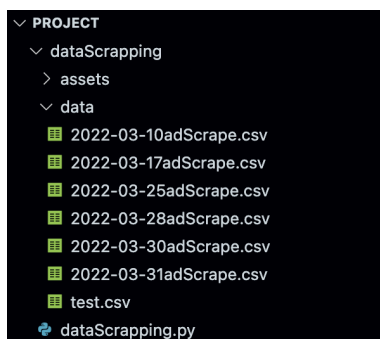
```

Код 5.4: CSV файлууд хадгалах

Дээрх код нь энгийн python програм файльтай харьцаж өөрт цуглуулсан өгөгдлөө хадгалж байна. Нийт өгөгдлийн хүснэгтийг энд ⁴ оруулав.

5.1.2 Ажлын байрны өгөгдөл

Zangia.mn ажлын байрны зарын вебсайтаас бакалаврын судалгааны ажлын үечилсэн төлөвлөгөөний дагуу 3 сарын 10-аас эхлэн өгөгдөл цуглуулсан билээ. Цуглуулж, өгөгдлийн сангийн зохиомжийн дагуу дараах утгуудын мэдээллийг CSV файлуудад хадгалж аваад байна. Доорх зурагт хамгийн



Зураг 5.2: Цуглуулсан өгөгдлийн файлууд

сүүлд буюу 3 сарын 31нд өгөгдлийн цуглуулга хийж 9174 өгөгдлийн excel хэлбэрт оруулсныг харж болж байна.

⁴<https://docs.google.com/spreadsheets/d/1rtATUKhUlleIKaWgFGvqiUWMipsrv-aCWZk-tYmzezU/edit?usp=sharing>

[illegible]

Зыпар 5.3: Data set

5.1.3 Өгөгдлийг нэгтгэх, цэвэрлэх

Ихэнх цуглуулсан өгөгдөл нь өгөгдлийн сангийн диаграмын дагуу амжилттай цуглуулсан бөгөөд дүн шинжилгээ хийх боломжтой өгөгдлүүдийг тусад нь хадгалж ашигласан болно.

Data scrape хийх явцад бүх өгөгдлийг *string* хэлбэрээр цуглуулсан бол өгөгдөлд дүн шинжилгээ

хийх явцад энэ нь тохиромжгүй тул цалинг *float* тохиролцох эсэхийг *boolean* болгож өөрчлөв.

```
1 def cleanSalary(salary) -> float:
2     if(isinstance(salary, str)):
3         return float(salary.replace(',', ''))
4     return None
5
6
7 def cleanDealable(deal) -> bool:
8     if(deal == ''):
9         return True
10    return False
```


11
12
13
14
15

```
def cleanLocation(location) -> str:

    if isinstance(location, str) and location != 'None':

        return location
```

Код 5.5: Өгөгдөл цэвэрлэх функц

	A	B	C	D	E	F	G	H	I	J	K
1		branch	employee	jobTitle	level	type	minSalary	maxSalary	isDealable	city	district
2		0 PR, олон нийтийн харилцаа	Надинбродерс Гадаад харилца	Мэргэжилтэн	Бүтэн цагийн		1200000	1500000	FALSE	Улаанбаатар хот	Баянзүрх дүүрэг
3		1 PR, олон нийтийн харилцаа	Оспринг Монгол Харилцааны ме	Мэргэжилтэн	Бүтэн цагийн		800000	1000000	TRUE	Улаанбаатар хот	Хан-Уул дүүрэг
4		2 PR, олон нийтийн харилцаа	Лэйдэр вэйт ХХ Олон нийттэй х	Мэргэжил хама	Бүтэн цагийн		1200000	1500000	TRUE	Улаанбаатар хот	Баянгол дүүрэг
5		3 PR, олон нийтийн харилцаа	Батбайгаль бэй Бүс хариуцсан	Мэргэжилтэн	Бүтэн цагийн		1000000	1200000	FALSE	Улаанбаатар хот	Баянзүрх дүүрэг
6		4 PR, олон нийтийн харилцаа	Бодь Интернэш ЭВЛҮҮЛЭГЧ, Э	Мэргэжилтэн	Бүтэн цагийн		1800000	2100000	FALSE	Улаанбаатар хот	
7		5 PR, олон нийтийн харилцаа	Тэнгэрийн Хишв График дизайн	Дунд шатны уд	Гэрээт/ Зөвлөх		1500000	1800000	FALSE	Улаанбаатар хот	Баянгол дүүрэг
8		6 PR, олон нийтийн харилцаа	Стандартформ Маркетингийн	Мэргэжилтэн	Бүтэн цагийн		800000	1000000	FALSE	Улаанбаатар хот	Баянзүрх дүүрэг
9		7 PR, олон нийтийн харилцаа	Гурван бухат ХО PR менежер	Мэргэжилтэн	Бүтэн цагийн		1500000	1800000	FALSE	Улаанбаатар хот	Сүхбаатар дүүрэг
10		8 PR, олон нийтийн харилцаа	Таван Богд Грүг PR менежер	Мэргэжилтэн	Бүтэн цагийн		1800000	2100000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг
11		9 PR, олон нийтийн харилцаа	Гурван бухат ХО Олон нийттэй х	Мэргэжилтэн	Бүтэн цагийн		1200000	1500000	TRUE	Улаанбаатар хот	Сүхбаатар дүүрэг
12		10 PR, олон нийтийн харилцаа	Жем Интернэш Олон нийт хари	Мэргэжилтэн	Бүтэн цагийн		1200000	1500000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг
13		11 PR, олон нийтийн харилцаа	Монполимет Гр ОЛОН НИЙТТЭ	Дунд шатны уд	Бүтэн цагийн		1500000	1800000	FALSE	Улаанбаатар хот	Сүхбаатар дүүрэг
14		12 PR, олон нийтийн харилцаа	Хүннү мебель Х Борлуулалтын	Дунд шатны уд	Бүтэн цагийн		1200000	1500000	FALSE	Улаанбаатар хот	Баянгол дүүрэг
15		13 PR, олон нийтийн харилцаа	Лавай Трейд ХО Худалдааны тө	Мэргэжил хама	Бүтэн цагийн		1500000	1800000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг
16		14 PR, олон нийтийн харилцаа	Лавай Трейд ХО Худалдааны тө	Мэргэжил хама	Бүтэн цагийн		1500000	1800000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг
17		15 PR, олон нийтийн харилцаа	Нутгийн буян гр ЭХ БЭЛТГЭГЧ	Мэргэжилтэн	Бүтэн цагийн		1800000	2100000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг
18		16 PR, олон нийтийн харилцаа	Грандмед эмнэл ХӨТӨЧ, БҮРТГ	Мэргэжилтэн	Бүтэн цагийн		800000	1000000	FALSE	Улаанбаатар хот	
19		17 PR, олон нийтийн харилцаа	Дрийм интерпр Маркетингийн	Дунд шатны уд	Бүтэн цагийн		1500000	1800000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг
20		18 PR, олон нийтийн харилцаа	Вертексмон ХХI Контент менеж	Дунд шатны уд	Бүтэн цагийн		1500000	1800000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг
21		19 PR, олон нийтийн харилцаа	Симатай Промс БАЙГУУЛЛАГА	Дунд шатны уд	Бүтэн цагийн		1000000	1200000	FALSE	Улаанбаатар хот	Баянзүрх дүүрэг
22		20 PR, олон нийтийн харилцаа	Лавай Трейд ХО Худалдааны тө	Мэргэжил хама	Бүтэн цагийн		1500000	1800000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг
23		21 PR, олон нийтийн харилцаа	Лавай Трейд ХО Худалдааны тө	Мэргэжил хама	Бүтэн цагийн		1500000	1800000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг
24		22 PR, олон нийтийн харилцаа	Оспринг Монгол Харилцагчийн	Мэргэжилтэн	Бүтэн цагийн		800000	1000000	TRUE	Улаанбаатар хот	Сүхбаатар дүүрэг
25		23 PR, олон нийтийн харилцаа	Эн Өү Ти Эс ХХ Утасны оператс	Мэргэжилтэн	Бүтэн цагийн		1200000	1500000	FALSE	Улаанбаатар хот	Баянгол дүүрэг
26		24 PR, олон нийтийн харилцаа	Эн Өү Ти Эс ХХ Утасны оператс	Мэргэжилтэн	Бүтэн цагийн		1200000	1500000	FALSE	Улаанбаатар хот	Баянгол дүүрэг
27		25 PR, олон нийтийн харилцаа	Ариг Банк	PR менежер	Мэргэжилтэн		1500000	1800000	TRUE	Улаанбаатар хот	
28		26 PR, олон нийтийн харилцаа	М Төсөл	Харилцагч хари	None		1200000	1500000	FALSE	Улаанбаатар хот	Сүхбаатар дүүрэг
29		27 PR, олон нийтийн харилцаа	Нью виндеу ХХI ХҮНИЙ НӨӨЦ,	Мэргэжилтэн	Бүтэн цагийн		1000000	1200000	FALSE	Улаанбаатар хот	
30		28 Авто засвар	ЮБИКСОЛЮШI Засварчин	None	Ээлжийн		1500000	1800000	FALSE	Төв аймаг	
31		29 Авто засвар	Намир ХХК	ТЕХНИКЧ	Мэргэжилтэн		1000000	1200000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг
32		30 Авто засвар	Сэрүүн сэлбэ Х Авто засварчин	Мэргэжилтэн	Бүтэн цагийн		1800000	2100000	TRUE	Улаанбаатар хот	Хан-Уул дүүрэг
33		31 Авто засвар	Грийн Групп	КУЗОВ ЗАСВАГ	Мэргэжилтэн		1800000	2100000	FALSE	Улаанбаатар хот	Сонгинохайрхан дүүрэг
34		32 Авто засвар	Автотерминал т Засварчин ажиг	None	Бүтэн цагийн		1200000	1500000	FALSE	Улаанбаатар хот	Хан-Уул дүүрэг

Зураг 5.4: Өгөгдлийг цэвэрлэж, өөрчлөлт хийсэн өгөгдөл

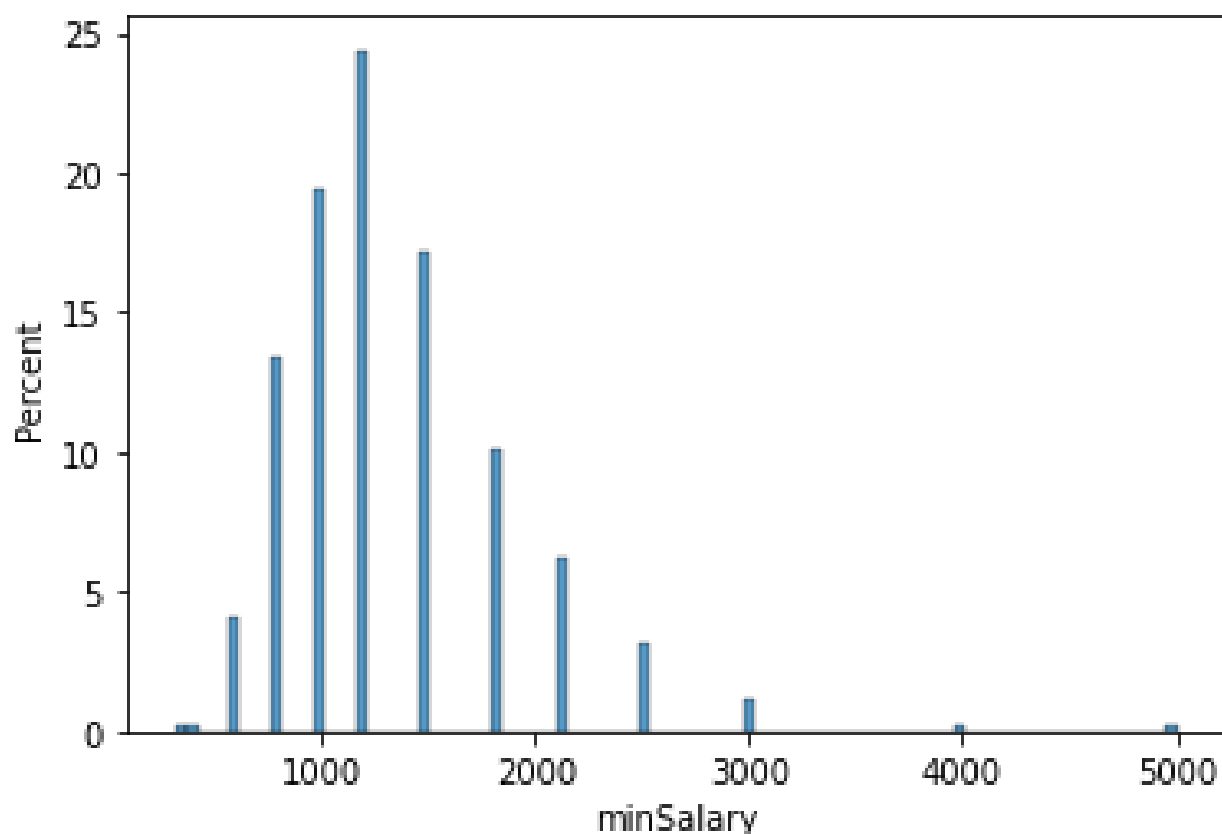
5.1.4 Өгөгдлийн статистик

Өөрчлөлт хийсэн өгөгдөл буюу нийт 8202 зарын өгөгдөл дээрээ дүн шинжилгээ хийж, цаашид үүнийгээ чатбот хөгжүүлэлт, асуулт хариултын загварт тусгахыг зорив.

Статистик - 1

Энэхүү графикт нийт цалингийн давхардсан утгууд нийт өгөгдлийн хэдэн хувийг эзэлж байгааг харуулж байна. Үүнээс харвал дундаж цалин 1 сая төгрөгөөс 1.5 сая төгрөгийн хооронд

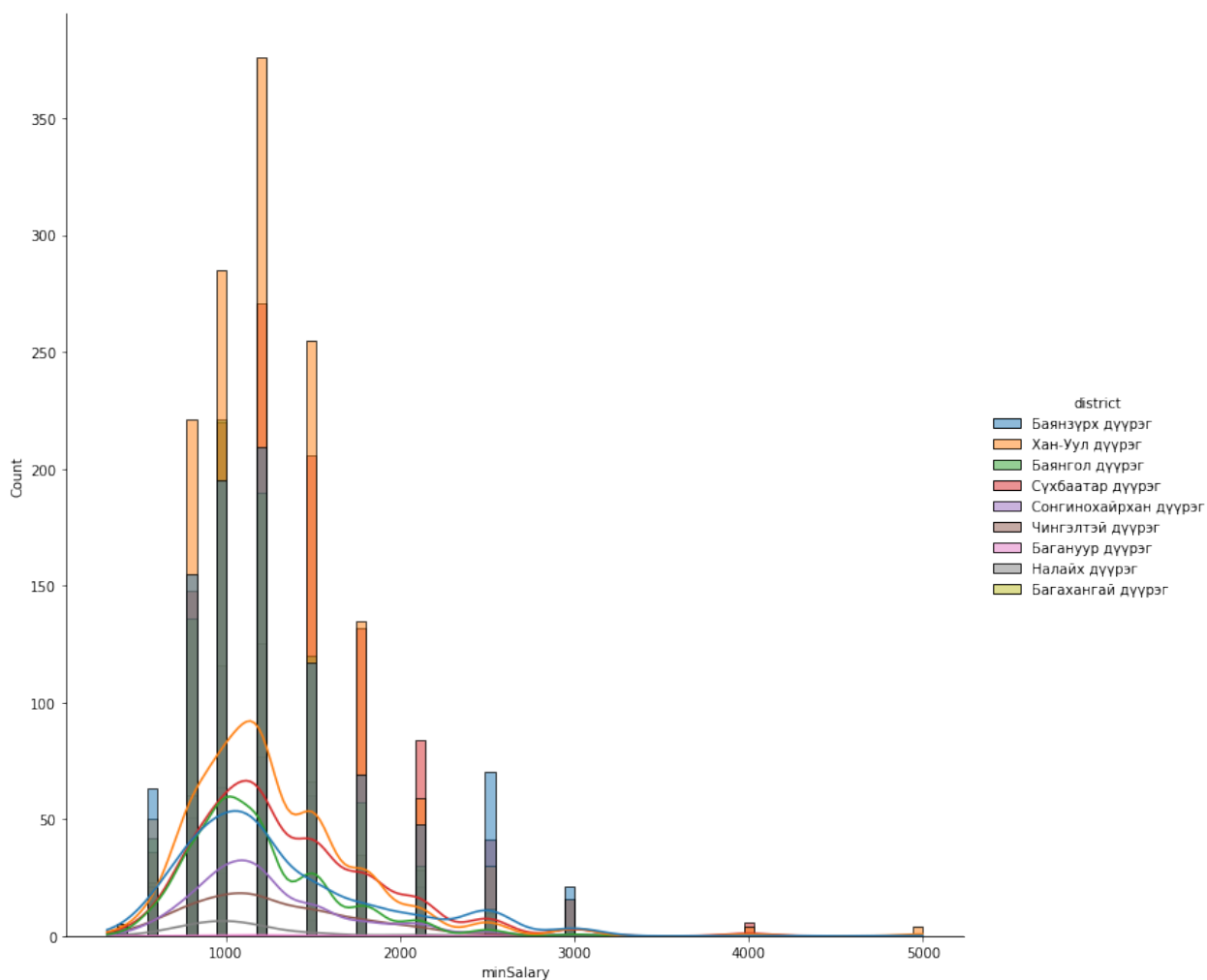
байгааг харж болж байна. Энэ нь нийт цалингийн хувийн 50-аас 55-ийг эзэлж байна.



Зураг 5.5: Өгөгдлийн статистик-1

Статистик - 2

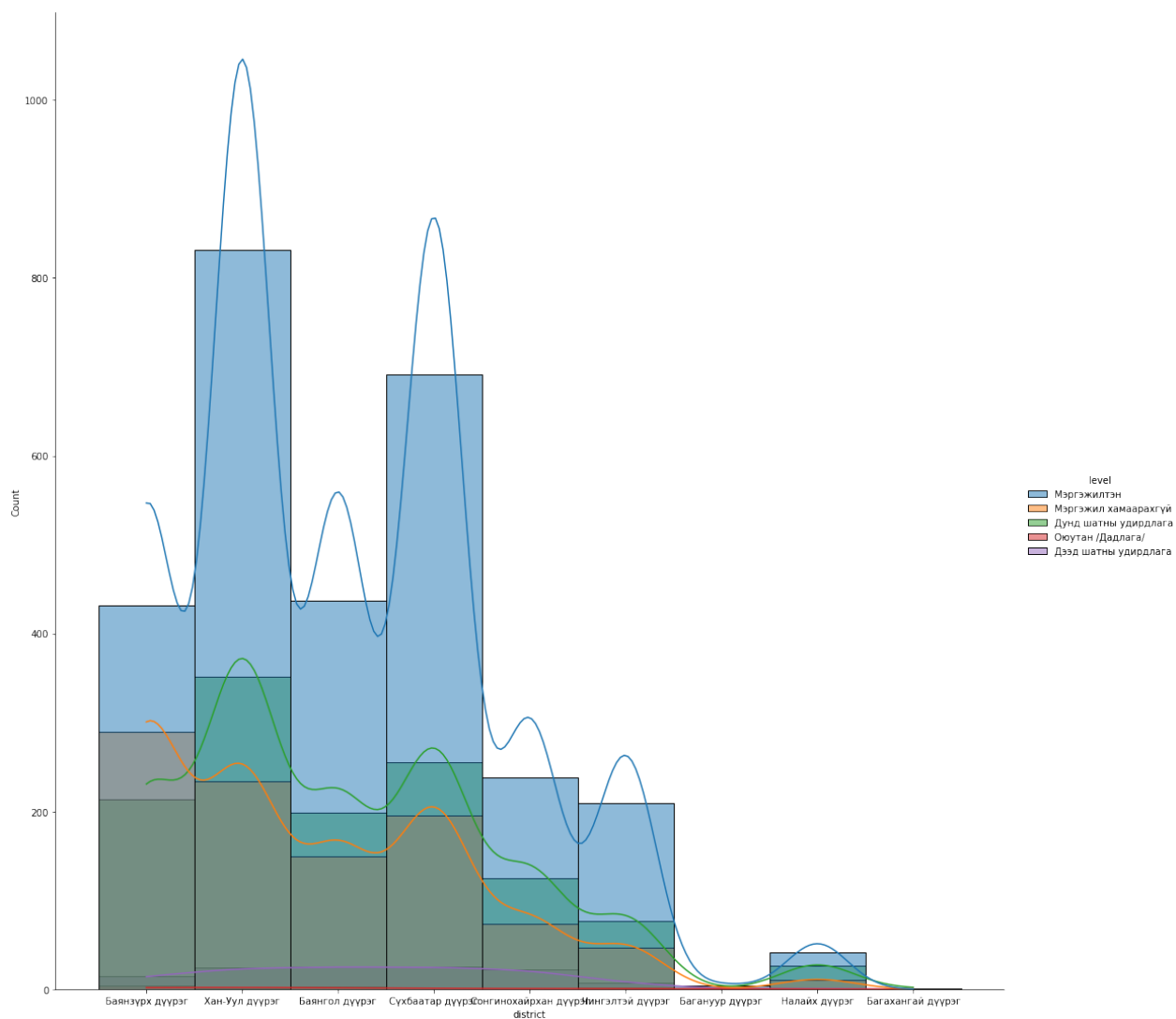
Доорх графикт нийт цалингийн давхардсан утгуудыг тоолж аль дүүрэгт хамгийн их байгааг өнгөөр нь дүрслэн харуулсан байна. Энэхүү графикаас харвал дундаж цалин буюу 1 сая төгрөгөөс 1.5 сая төгрөгөөр цалинжуулах олон ажлын байр санал болгож байгаа дүүрэг нь Хан-Уул болон Сүхбаатар дүүрэг байна. Ажлын байрны төвлөрөл болон их хотын бүтээгдэхүүн үйлдвэрлэлийн цэгийг Хан-Уул, Сүхбаатар дүүрэг гэж дүгнэж болохоор байна.



Зураг 5.6: Өгөгдлийн статистик-2

Статистик - 3

Өмнөх графикийн дүгнэлтийн адилаар хамгийн их ажлын санал болгож буй дүүрэг нь Хан-Уул, Сүхбаатар дүүрэг байх бөгөөд ажлын байрны шаардах түвшинг хамтад нь харуулсан байна. Үүнээс үзвэл, мэргэжилтэн болон дунд шатны удирдлагын орон тоо эрэлттэй байна гэж үзэж болно. Харин дадлагажигч болон дээд шатны удирдлагын эрэлт харьцангуй бага байгааг харж болж байна.



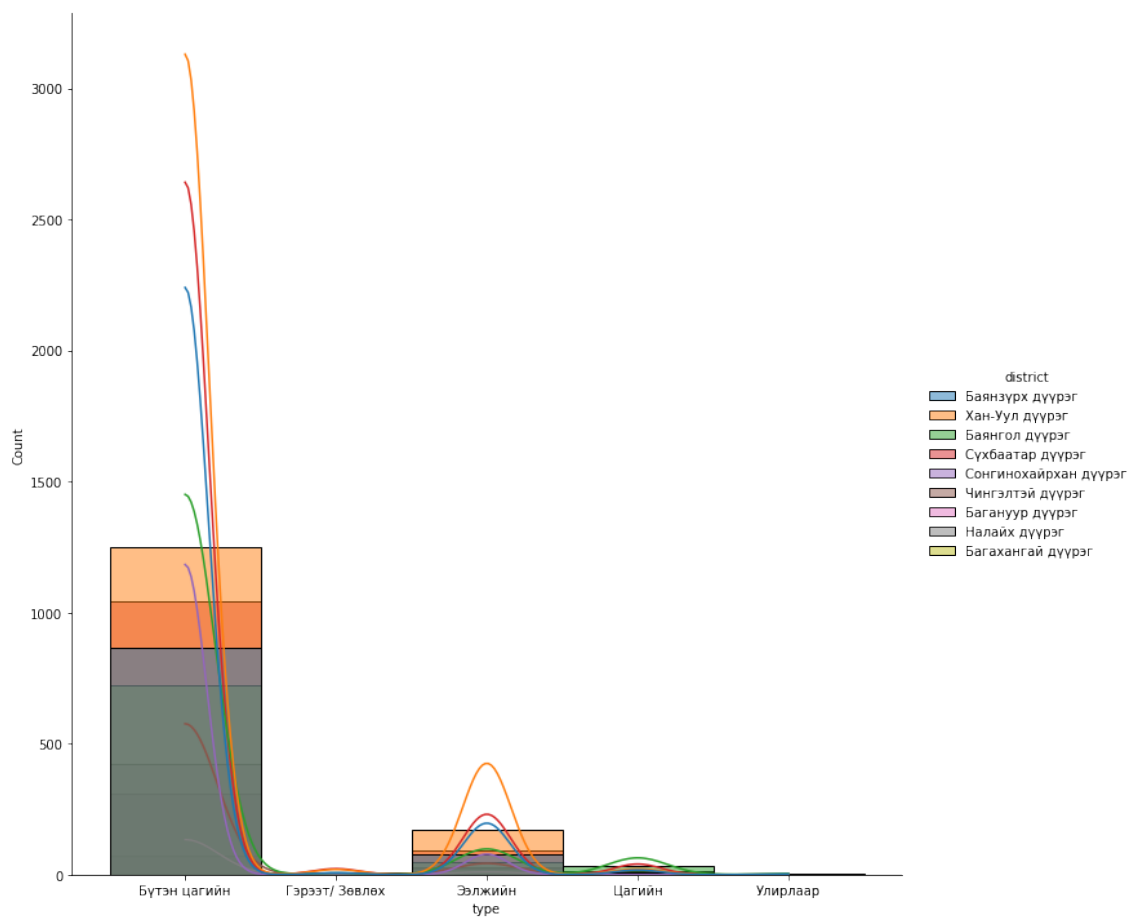
Зураг 5.7: Өгөгдлийн статистик-3

Статистик - 4

Энэхүү графикт дүүргүүд харгалзан ямар төрлийн цагийн хуваарьтай ажил санал болгож байгаа болон тэдгээрийн тоотой нь харьцуулан дүрсэлжээ. Эндээс бүтэн цагийн ажилтан болон ээлжийн төрлийн ажлын байр ихэнх хувийг эзэлж байгааг харлаа.

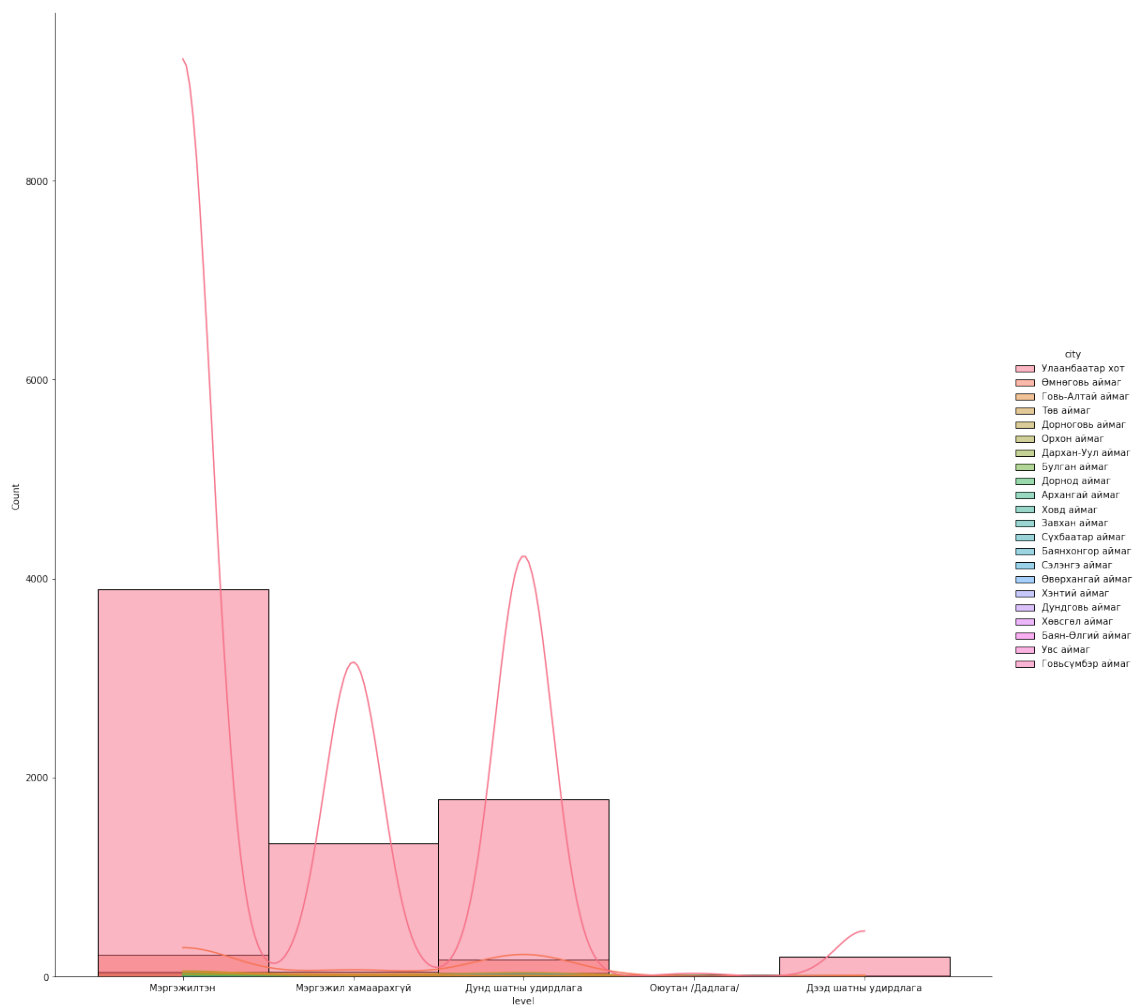
Статистик - 5

Доорх графикаас ажил олгогчид ямар төрлийн цагийн хуваарьтай, хаана ажил санал болгож байгааг 21 аймгаар бүсчлэн өнгөөр илэрхийлсэн байна. Үүнээс дүгнэвэл, Улаанбаатар хотод



Зураг 5.8: Өгөгдлийн статистик-4

нягтаршил маш өндөр байгаа бөгөөд ажлын байрны эрэлт аймгуудтай харьцуулахад маш өндөр байна.



Зураг 5.9: Өгөгдлийн статистик-5

Дүгнэлт

Бакалаврын судалгааны ажлаар “Ажил олгогчдын өгөгдлийн анализ систем дээр суурилсан чатбот” сэдвийн дагуу хөгжүүлэлтийг эхлүүлсэн бөгөөд уг судалгааны ажилд холбогдох онолын судалгаа, ашиглаж буй технологи, түүнийг илүү онолын мэдлэг болон системийн хөгжүүлэлтийн хэсгээс дэлгэрэнгүй тайлбарласан болно.

Bibliography

- [1] Чатбот системийн тухай
<https://www.engati.com/blog/types-of-chatbots-and-their-applications>
- [2] Өгүүлбэр хувиргалтын арга зүй
<https://www.sbert.net/docs/quickstart.html>
- [3] Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks
<https://arxiv.org/abs/1908.10084>
- [4] Use case diagram
https://app.diagrams.net/#G1jhom3sc_holt-X9XLALtQja_Gl_Eykhj
- [5] Business Process Model Notation 2.0 диаграмм
<https://cawemo.com/diagrams/ea037ec0-c1c5-4ab6-8262-521657472803--bpmn-2-0?v=960,418,1>
- [6] Өгөгдлийн сангийн диаграмм
<https://dbdiagram.io/d/6249fb7cd043196e39e87451>

А. ҮЕЧИЛСЭН ТӨЛӨВЛӨГӨӨ

Батлаа.

МКУТ-ийн эрхлэгч:...../док. проф. Н.Оюун-Эрдэнэ/

2022 оны 02 сарын 11

Монгол нэр Ажил олгогчдын өгөгдлийн анализ систем дээр суурилсан чат бот

Англи нэр Chat bot based on system analysis of employers' data

Сэдэвт бакалаврын судалгааны ажлын 7 хоногийн үечилсэн төлөвлөгөө

Хугацаа: 2022.02.07-оос 2022.05.06 хүртэл 13 долоо хоног

№	Хийх ажил	Долоо хоног													14 Жинхэнэ хамгаалалт	Тайлбар
1	Онолын судалгаа															
	Scrapper tool															
	Bot tool															
2	Өгөгдөл цуглуулалт															
	Цуглуулах код бичих															
	Өгөгдлийг бааруулах															
3	Системийн шаардлага тодорхойлох															
	Хэрэглэгчийн шаардлага тодорхойлох															
5	Системийн зохиомж															
	Өгөгдлийн сэнтийн зохиомж															
	Чат бот хөгжүүлэлт															
6	Хэрэгжүүлэлт															
	Өгөгдөлд анализ хийх хайгуулах															
7	Бичиг баримт															
	Тайлан боловсруулах															

Тайлбар: Төслийг гэрээжүүлэх төлөвлөгөөг 7 хоногийн дотоод хамгаалалтаар хийж төв гэрээг бүрдэж гүйцэтгэнэ. Хийх ажил дэд хэсгийн байдал үз ажилд зарцуулах хугацааг хувиар тэмдэглэж байна. Ажлын эхлэх өдөр 2022.02.07-оос 2022.05.06 хүртэл 13 долоо хоног.

Зөвшөөрсөн: Удирдагч багш

Боловруулсан: Оюутин

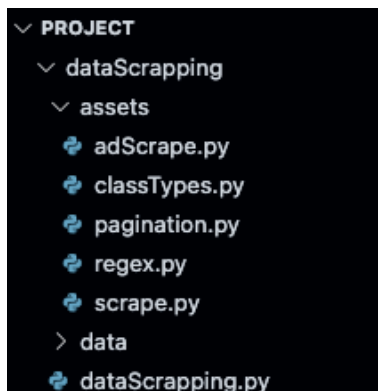
Мэдээллийн технологи А. Сайнболбоо/ Оюутны ID: 1851num1762 Холбогдох утас: 91990388

Зураг А.1: Бакалаврын судалгааны ажлын үечилсэн төлөвлөгөө

В. КОДЫН ХЭРЭГЖҮҮЛЭЛТ

В.1 Өгөгдөл цуглуулалт

Өгөгдөл цуглуулах програм нь дараах бүтэцтэй байх бөгөөд assets доторх кодууд нь үндсэн кодыг ажлуулахад туслах функцууд байна.



Зураг В.1: Фолдерийн бүтэц

В.1.1 Үндсэн өгөгдлийг цуглуулах эх код

```
1 from datetime import date
2 import time
3 from assets.classTypes import Category
4 from assets.scrape import UseBeautifulSoup as useScrape
5 from assets.adScrape import advertisementScrape as useAdScrape
6 from assets.pagination import createLinkList as createLinkList
7
8 start_time = time.time()
9 initialUrl = 'https://www.zangia.mn/'
10 today = str(date.today())
11 # all categories set
12 categorySet = set()
13 # all advertisement's link set
14 adUrlDict = {}
15 # all ads object set
16 adsSet = set()
17
18 # scrape initial links
19 soup = useScrape(initialUrl)
20 navigatorList = soup.find_all('div', class_='filter')
21 for navigator in navigatorList:
22     if navigator.find('h3').text.strip() != 'Salbar, mergejl':
23         continue
24     # ALL CATEGORY LINKS
25     categoryList = navigator.find_all('div')
```

```

26
27 for categoryItem in categoryList:
28     categories = categoryItem.find('a')
29     url = initialUrl + categories['href']
30     tempCategory = Category(url, categories.text, '')
31     soup = useScape(url)
32     subCategory = soup.find('div', class_='pros')
33     # ALL SUBCATEGORY LINKS
34     subCategoryList = subCategory.find_all('a')
35     for subCategoryItem in subCategoryList:
36         subCategoryUrl = initialUrl + subCategoryItem['href']
37         tempSubCategory = Category(
38             subCategoryUrl, subCategoryItem.text, tempCategory.name)
39         categorySet.add(tempSubCategory)
40
41 for categoryItem in categorySet:
42     if categoryItem.parentId == '':
43         continue
44     soup = useScape(categoryItem.url)
45     hasPagination = soup.find('div', class_='page-link')
46     pagesUrl = []
47     if hasPagination != None:
48         pagesUrl = createLinkList(hasPagination, categoryItem.url)
49     else:
50         pagesUrl.append(categoryItem.url)
51     for pageUrl in pagesUrl:
52         soup = useScape(pageUrl)
53         ads = soup.find_all('div', class_='ad')
54         # CREATE UNIQUE AD DICTIONARY
55         for ad in ads:
56             adUrl = initialUrl+ad.find('a', class_=None)['href']
57             adUrlDict[adUrl] = categoryItem
58     pagesUrl.clear()
59
60 file = open(today+'adScrape.csv', 'w', encoding='utf-8')
61 file.write('Parent Category Name' + '\t' +
62           'Category Name ' + '\t' +
63           'Link' + '\t' +
64           'Employee Company' + '\t' +
65           'Title' + '\t' +
66           'Roles' + '\t' +
67           'Requirements' + '\t' +
68           'Additional Info' + '\t' +
69           'City/Province' + '\t' +
70           'District' + '\t' +
71           'Level' + '\t' +
72           'Type' + '\t' +
73           'Min Salary' + '\t' +
74           'Max Salary' + '\t' +
75           'Is Dealable' + '\t' +
76           'Address' + '\t' +
77           'Phone' + '\t' +

```

```

78         'Fax' + '\t' +
79         'Ad Added Date' + '\n')
80
81 for adUrl in adUrlDict:
82     print(adUrl)
83     try:
84         tempAdItem = useAdScrape(adUrl)
85         tempAdItem.setCategory(adUrlDict[adUrl])
86         file.write(
87             tempAdItem.category.parentId+'\t' +
88             tempAdItem.category.name+'\t' +
89             tempAdItem.url+'\t' +
90             tempAdItem.company+'\t' +
91             tempAdItem.title+'\t' +
92             tempAdItem.roles+'\t' +
93             tempAdItem.requirements+'\t' +
94             tempAdItem.additionalInfo+'\t' +
95             tempAdItem.city+'\t' +
96             tempAdItem.district+'\t' +
97             tempAdItem.level+'\t' +
98             tempAdItem.type+'\t' +
99             tempAdItem.minSalary+'\t' +
100            tempAdItem.maxSalary+'\t' +
101            tempAdItem.isDealable+'\t' +
102            tempAdItem.address+'\t' +
103            tempAdItem.phoneNumber+'\t' +
104            tempAdItem.fax+'\t' +
105            tempAdItem.adAddedDate+'\n')
106         del tempAdItem
107     except:
108         print('Ad writing error')
109 file.close()
110 print("--- %s seconds ---" % (time.time() - start_time))

```

Код В.1: Бүх өгөгдлийг цуглуулах - dataScraping.py

В.1.2 Нэг зарын шаардлагатай бүх мэдээллийг цуглуулах код

```

1 import re
2 from .classTypes import Advertisement
3 from .scrape import UseBeautifulSoup as useScrape
4
5
6 def listScrapper(sections, key) -> str:
7     content = []
8     for section in sections:
9         subTitle = section.find('h2', class_=None).text
10        if key != subTitle:
11            continue
12        div = section.find('div', class_=None)
13        children = div.next_element
14

```

```

15         while(children != None):
16             try:
17                 content.append(textStrip(children.text))
18                 children = children.next_sibling
19                 continue
20             except:
21                 print('An error occurred')
22                 children = children.next_sibling
23             content = [s for s in filter(listFunc, content)]
24         if not content:
25             return ''
26         return ' '.join(content)
27
28
29 def textStrip(text) -> str:
30     pattern = re.compile('[\r\n\xa0\t ]+', re.MULTILINE | re.IGNORECASE)
31     return pattern.sub(' ', text.strip())
32
33
34 def listFunc(e):
35     return len(e) != 0
36
37
38 def singleItemScrapper(sections, key, subKey) -> str:
39     for section in sections:
40         subTitle = section.find('h2', class_=None).text
41         if key != subTitle:
42             continue
43         div = section.find_all('div', class_=None)
44         for item in div:
45             if item.next_element.text == subKey:
46                 return textStrip(item.find('span').text)
47     return 'None'
48
49
50 def salaryScrapper(salary):
51     isDealable = ''
52     k = re.split(r'[^d,]+', salary, 2, re.IGNORECASE)
53     if len(k) < 2:
54         [a] = k[0:1]
55         return a, a
56     [a, b] = k[0:2]
57     if len(k) > 2:
58         isDealable = ' '
59     return a, b, isDealable
60
61
62 def locationScrapper(location):
63     city = ''
64     district = ''
65     k = location.split(',')

```

```

66     if len(k) < 2:
67         city = k[0]
68         return city, district
69     [city, district] = k[0:2]
70     return city, district
71
72
73 def advertisementScrape(url) -> Advertisement:
74     soup = useScrape(url)
75     advertisement = Advertisement(url, soup.find('h3').text.strip())
76     companyTitle = soup.find('div', class_='nlp').find('td')
77     for item in companyTitle:
78         try:
79             if item.name == None:
80                 advertisement.company = textStrip(item.text)
81         except:
82             print('Company name scrape error')
83     # advertisement.company = textStrip(company)
84
85     # all items
86     sections = soup.find_all('div', class_='section')
87     advertisement.roles = listScrapper(
88         sections, 'Guitsetgeh undsen uurg'')
89     advertisement.requirements = listScrapper(
90         sections, 'Ajliin bairnii shaardlaga')
91     advertisement.additionalInfo = listScrapper(
92         sections, 'Nemelt medeelel')
93     advertisement.level = singleItemScrapper(sections, 'Busad', 'Tuvshin
94         ')
95     advertisement.type = singleItemScrapper(sections, 'Busad', 'Turul')
96     minSalary, maxSalary, isDeable = salaryScrapper(
97         singleItemScrapper(sections, 'Busad', 'Tsalin'))
98     city, district = locationScrapper(
99         singleItemScrapper(sections, 'Busad', 'Bairshil'))
100     advertisement.minSalary = minSalary
101     advertisement.maxSalary = maxSalary
102     advertisement.isDeable = isDeable
103     advertisement.city = city
104     advertisement.district = district
105     advertisement.address = singleItemScrapper(sections, '
106         ', ' ')
107     advertisement.phoneNumber = singleItemScrapper(
108         sections, 'Holboo barih', 'Utas')
109     advertisement.fax = singleItemScrapper(
110         sections, 'Holboo barih', 'Fax')
111     advertisement.adAddedDate = singleItemScrapper(
112         sections, 'Zariin hugatsaa', 'Zar niitelsen ognoo')
113     print(advertisement.additionalInfo)
114     print('SINGLE AD SCRAPING DONE!!!', url)
115
116     return advertisement

```

Код В.2: Нэг зарын өгөгдлийг цуглуулах - adScrape.py

B.1.3 Цуглуулах өгөгдлийн төрөл

```
1 class Category:
2     url = ''
3     name = ''
4     parentId = ''
5
6     def __init__(self, url, name, parentId='None') -> None:
7         self.url = url
8         self.name = name
9         self.parentId = parentId
10
11     def getUrl(self) -> str:
12         return self.url
13
14 class Advertisement:
15     category = Category
16     url = ''
17     company = ''
18     title = ''
19     # ListInfo
20     roles = ''
21     requirements = ''
22     additionalInfo = ''
23     # OtherInfo
24     city = ''
25     district = ''
26     level = ''
27     type = ''
28     minSalary = ''
29     maxSalary = ''
30     isDealable = ''
31     # ContactInfo
32     address = ''
33     phoneNumber = ''
34     fax = ''
35     adAddedDate = ''
36
37     def __init__(self, url, title) -> None:
38         self.url = url
39         self.title = title
40
41     def setCategory(self, category) -> None:
42         self.category = category
```

Код B.3: Өгөгдлийн төрөл - classTypes.py

B.1.4 BeautifulSoup scraper

```
1 from bs4 import BeautifulSoup
2 import requests
3 from urllib.error import HTTPError
```

```

4
5
6 def UseBeautifulSoup(url):
7     try:
8         response = requests.get(url)
9         response.raise_for_status()
10    except HTTPError as error:
11        print(error)
12    soup = BeautifulSoup(response.text, 'html.parser')
13    return soup

```

Код В.4: Scrape хийх функц - scrape.py

В.2 Өгөгдөл нэгтгэх, цэвэрлэх

```

1 from random import random
2 import requests
3 import pandas as pd
4 import numpy as np
5 import csv
6 import bs4 as BeautifulSoup
7 import random
8
9 with open('/Users/zolboo/Desktop/Bachelor/employmentAnalysis/project/
    latestData.tsv', 'r', encoding='utf-8') as file:
10     app_lines = file.read().split('\n')
11
12 df = pd.read_csv(
13     '/Users/zolboo/Desktop/Bachelor/employmentAnalysis/project/
        latestData.tsv', sep='\t')
14
15
16 def cleanSalary(salary) -> float:
17     if(isinstance(salary, str)):
18         return float(salary.replace(',', ''))
19     return None
20
21
22 def cleanDeable(deal) -> bool:
23     if(deal == ''):
24         return True
25     return False
26
27
28 def cleanLocation(location) -> str:
29     if isinstance(location, str) and location != 'None':
30         return location
31     return None
32
33
34 def normalizeDataSet(data_set):
35     ret = pd.DataFrame(columns=['branch', 'employee', 'jobTitle', '

```



```

36         level',
37         'type', 'minSalary', 'maxSalary', 'isDealable',
38         'city', 'district'])
39     for index, row in data_set.iterrows():
40         minSalary = cleanSalary(row['Min Salary'])
41         maxSalary = cleanSalary(row['Max Salary'])
42         location = cleanLocation(row['City/Province'])
43         print(type(row['District']), row['District'])
44         if minSalary is None and maxSalary is None and location is None
45             :
46             continue
47         dealable = cleanDealable(row['Is Dealable'])
48         ret = ret.append({'branch': row['Category Name'],
49                         'employee': row['Employee Company'],
50                         'jobTitle': row['Title'],
51                         'level': row['Level'],
52                         'level': row['Type'],
53                         'minSalary': minSalary,
54                         'maxSalary': maxSalary,
55                         'isDealable': dealable,
56                         'city': location,
57                         'district': row['District']}
58                         , ignore_index=True)
59     return ret
60
61 data_set = normalizeDataSet(df)
62 data_set.to_csv('data.csv')

```

Код В.5: Өгөгдөл цэвэрлэх - dataClean.py