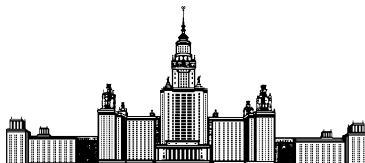


Московский государственный университет имени М. В. Ломоносова



Факультет Вычислительной Математики и Кибернетики
Кафедра Математических Методов Прогнозирования

ВЫПУСКНАЯ КВАЛИФИКАЦИОННАЯ РАБОТА

**«Методы обучения без учителя для выделения поляризации
в новостных потоках»
«Unsupervised polarization detection methods in newsflows»**

Выполнил:

студент 4 курса 417 группы

Висков Василий Алексеевич

Научный руководитель:

д.ф-м.н., и.о. заведующего кафедры

Воронцов Константин Вячеславович

Москва, 2022

1	Введение	3
1.1	Основные определения	4
1.2	Исходные данные	5
1.3	Обзор области	7
1.4	Цель работы	9
1.5	Постановка задачи машинного обучения	9
2	Решение	12
2.1	Кластеризация	12
2.1.1	Алгоритмы, основанные на анализе плотностных характеристик выборки	12
2.1.2	Алгоритмы, явно параметризуемые количеством кластеров	14
2.1.3	Выводы	17
2.2	Векторное представление	19
2.3	Снижение размерности	20
2.4	Разметка данных	21
3	Вычислительный эксперимент	22
3.1	Результаты эксперимента	22
3.2	Выводы	25
4	Исследование абляции	25
4.1	Результаты эксперимента	26
4.2	Выводы	27
5	Заключение	27
6	Приложение	30
6.1	Инструкция для ассессоров для Яндекс.Толока	30

Аннотация

В данной работе исследуется задача на стыке **Polarization Detection** и **Stance Detection**, предлагается алгоритм разбиения выборок документов, связанных одним общим событием, на непересекающиеся подмножества, описывающие мнения, с нефиксированным их количеством методами обучения без учителя и предлагается способ формирования поляризованных размеченных выборок.

Важные общественные события, связанные с политикой, спортом или другими насущными темами, освещаются рядом журналистских служб, и зачастую публикуемые новости описывают произошедшее с разных позиций, ссылаясь либо на одну из противоборствующих сторон, либо вообще субъективно оценивая ситуацию. В рамках новостной полемики вокруг некоторого события могут возникать различные мнения, зачастую выраженные не столько используемой лексикой, сколько эмоциональной окраской и степенью эксцентричности, с которой смысл текста раскрывается. Причем этих мнений может быть произвольное количество в зависимости от степени общественного резонанса в рамках новостного фона.

Рассмотрим несколько примеров, демонстрирующих поляризованность новостного потока вокруг некоторого события.

1. Тема: «Разрушение Израилем офисного здания в секторе Газа»

- Полюс 1: «Взгляд одной из телевизионных компаний, базирующейся в этом здании»
 - «В «Аль-Джазире» прокомментировали уничтожение Израилем своего офиса в Газе. Там подчеркнули, что эти действия являются «варварским актом» который направлен на то, чтобы помешать «говорить правду».
- Полюс 2: «Взгляд израильской стороны»
 - «Израиль нанес удары по второму высотному зданию в Газе, сообщил источник. Как позднее пояснили в Армии обороны Израиля, в этом доме находилось военное оборудование разведки ХАМАС, а офисы медиакомпаний движение использовало как живой щит».
- Полюс 3: «Нейтральный взгляд/констатация факт»
 - «Израиль разрушил здание в Газе с офисами международных СМИ. Авиаудар, нанесенный Израилем по сектору Газа, привел к обрушению 15-этажного здания. Там, в частности, располагались офисы международных СМИ. Речь идет о высотном здании «аль-Джала».

2. Тема: «Инцидент с пермскими подростками и таксистом»

- Полюс 1: «Произвол/негатив»
 - «Пермские подростки изрезали водителя такси».
- Полюс 2: «Замалчивание»
 - «В Перми студенты колледжа ПГНИУ катались в машине с раненым таксистом».
- Полюс 3: «Констатация факта»
 - «Трое студентов колледжа в Перми стали фигурантами дела о нападении с ножом на таксиста».

В данной работе будут рассмотрены способы построения кластеризации для произвольной выборки документов из новостного потока, объединенных описанием одного события, где кластер, или полюс, характеризует поляризованное или нейтральное мнение, а также шумовое подмножество документов не о событии выборки. Мнение должно быть связано с одним и только одним событием и может характеризоваться не только конкретной позицией о каком-то вопросе, но и степенью эксцентричности, с которой эта позиция раскрывается. Также будет построена совокупность размеченных валидационных выборок на русском языке, состоящих из подмножеств документов, объединенных одним событием и включающих тексты, описывающие в совокупности от 1 до 10 различных мнений.

1.1

Разбиение множества – это представление его в виде объединения произвольного количества попарно непересекающихся непустых подмножеств.

Темой будем называть множество документов, большинство из которых описывает одно и то же событие. Это событие будем называть **событием** темы.

Поляризованное мнение (мнение), или **полюс поляризации** – эмоционально, синтаксически или семантически однородное подмножество исходной темы, задающее некоторую позицию вокруг события. Будем считать, что нейтральное мнение также является полюсом поляризации.

Документы темы, не относящиеся к его событию, будем называть **шумовыми**. Совокупность мнений и шумовых документов образуют разбиение темы.

Привилегированная информация — дополнительная информация об объектах выборки, доступная только на этапе обучения. Таковой в данной работе выступает разметка подмножества генеральной совокупности тем, коими являются и предлагаемые в дальнейшем к разметке наборы документов. Эти темы будут использоваться исключительно для измерения качества разбиения их на полюсы поляризации.

Тонально окрашенными будем называть слова или словосочетания, выражающие эмоциональное отношение автора высказывания к некоторому объекту, выраженном в тексте.

Тональностью слова или словосочетания будем называть его степень тональной окраски и определять целочисленной величиной.

Именованной сущность называют слово или словосочетание, четко идентифицирующее один элемент из набора других элементов, имеющих аналогичные свойства. Примерами именованных сущностей могут быть ФИО, названия мест и компаний, т.д.

Эмбединг (от англ. *embedding*) – вещественнозначный вектор некоторой размерности, выступающий признаковым описанием объекта (в данной работе – текста).

Матрицей объектов-признаков выборки – мощности $\#$ называется вещественная матрица $\in \mathbb{R}^{\# \times n}$, состоящая по строкам из векторных представлений $\in \mathbb{R}^n$ объектов выборки.

1.2

Предоставленными данными служит совокупность документов из новостного потока X . Внешним алгоритмом по отношению к этой работе они разбиты на совокупность тем $\{t_j\}_{j=1}^{\%}$, где $\% = 90$. Каждая из тем $t_j = \{s_1^j, \dots, s_{n_j}^j\}$, где s_i^j – один из признаков документа. Темы описывают события из следующих новостных рубрик:

- происшествия (кол-во: 42);
- наука и техника (кол-во: 12);
- политика (кол-во: 27);
- культура (кол-во: 3);

^ ñèèíáúá ñòðóèòóóú (éíè-âí: 3); ^ ýéíííèèà/òèíáíñú (éíè-âí: 3);

Ñíáúòèyìè ìíãòò ñéóæèòü, íàíðèíáð, ¼ãñòðíííú NASA áíáðáúá íáíàðóæèèè ðáíòãáííñèèá èó÷è ìò óðàíàç èèè ¼Áàéááí ìèèàçàèñý ìò áñòðá÷è ñ Çàèáíñèè ìáðáá ðàçãíáíðí ñ Íóèèííç.

Đèñ. 1: Èññèääíááíèá àáííúó

Ñðááíýý äèèà áíèóíáíòà ñ ó÷àòí çàãíèíèè è òàèà ñíñòàáèýáò ìðèíáðíí 1388 ñèíáíèíá èèè 188 ñèíá. Èçããñòíí¹, ÷òí ñðááíýý äèèà íóáèèèèèèè á ñàòè Twitter ñíñòàáèýáò 28 ñèíáíèíá, ìàèñèíàèüíáý äèèà ìðèíáðíí 140 ñèíáíèíá, ÷òí íá èääò íè á èàèíá ñðááíáíèá ñ èìàðüèèèñý àáííúè.

Èíèè÷ãñòáí áíèóíáíòá á òáíàð ñèèüíí ñíáúáíí á ñòíðíó íááíèüøèò ááèè÷èí, äèý íááíèüøèò áúáíðíè ñíñòðíáíèá ìðáàèèèíè èèãñòðèçàòèè ìæàò áúòü ñèíæíé çààà÷áé.

Ìðèçíàèèè áíèóíáíòá áúñòóíàðò, ìíèí çàãíèíèè è òàèà áíèóíáíòà, è ðàçóèüòàòü íáðááíðèè ñííèíèóíííñòüð ìíááèáé ìíèíáíèý áñòáñòááíííáí ýçúèà, áíáø-íèò ìí ìòííðáíèð è ýóíé ðááíðá, à èíáíí: ðàçíáòèè òàèà áíèóíáíòà íà èíáííáí-

¹<https://smk.co/the-average-tweet-length-is-28-characters-long-and-other-interesting-facts/>

íúá nõúííñòè ñ ìðíñòààèèáíèàì èàòááíðèè ¼èè÷ííñòù¿, ¼ääíãðàòè÷áñèíà ìíèíæáíèà¿ èèè ¼íðääíèçàòèý¿, ðàçíàòèà òííàèüííñòè äèý èìáííàáííúò nõúííñòàé è ñíòèèèüíí-ääííãðàòè÷áñèèà ìðèçíàèè, ñâyçàííúá ñ òíðìèðííààíèàì òäèääíé àóäèòíðèè äèý èàæáííà èç áíèòíáíòíá. Ìíñèääíèà ìðèçíàèè ìæíí ðàçáèòù íà äðóííú:

- ^ ááðíýòííñòù òíáí, ÷òí òäèääíé àóäèòíðèèé áíèòíáíòà ýäèýàòñý íáèíòíðàý áíç-ðàñòíàý äðóííà (áúääèýðòñý 4 äðóííú: ìò 0 áí 21, ìò 22 áí 39, ìò 40 áí 59, ìò 60 è ñòàððá);
- ^ ááðíýòííñòù òíáí, ÷òí òäèääíé àóäèòíðèèé áíèòíáíòà ýäèýàòñý äðóííà èðááé, èìàðùèò ìíðääèáííóð àèääáíè÷áñèóð ñòáíáíú èèè íáðàçííàíèà (áúääèýðòñý 5 äðóííú: áíèòíðà íáóè, èáíàèääòù íáóè, áúñøáá, ñðááíáá íáúáá è ñðááíáá ñíàòèèèüííà íáðàçííàíèà);
- ^ ááðíýòííñòù òíáí, ÷òí òäèääíé àóäèòíðèèé áíèòíáíòà ýäèýàòñý äðóííà èðááé ñ ìíðääèáííúì áæáíàñý÷íúì áíñòàòèíí, áúðàæáííí á íáèíòíðùò óñèíáíúò áäèíèòàò (áúääèýðòñý 4 äðóííú: áíñòàòíè ìò 0 áí 50, ìò 50 áí 100, ìò 100 áí 150 è ìò 150 áí 200);
- ^ ááðíýòííñòù òíáí, ÷òí òäèääíé àóäèòíðèèé ýäèýðòñý æáíúèíú èèè íóæ÷èíú;

Òííàèüííñòù ýäèýàòñý ìíðýäèíáúì ìðèçíàèè íà áííáíá òäèúò ÷èñáè, ááá 0 ñíò-áàòñòáóàò íáèòðàèüííó ìòííøáíèð, ìíáóèü çíà÷áíèý ìòíáðàæáàò ñòáíáíú òííàèüííé ìèðàðáíííñòè äèý èàæáíé èç èìáííàáííúò nõúííñòàé, à çíàè ìòðèòàòäèüííà èèè ìíèíæèòäèüííà ìòííøáíèà è íáúáèòó ìðèáýçèè.

Éíèè÷áñòáí íáííèé á èàæáíé èç òáí íáèçááñòíí, èàè ñðááè èìàðùèòñý ááííúò, òàè è áí áñáé ááíáðàèüííé ñíííèòíííñòè òáí, ìíýòííó ðáðàòù çàää÷ó íàòíààè, èí-òíðùá èàèè-èèáí íáðàçíí çàááèñòáóðò ðàçíàòèó, ò.á íàòíààè íáó÷áíèý ñ ó÷èòäèáí èèè íàòíààè ÷áñòè÷ííá íáó÷áíèý, íà ìðááñòàèèýàòñý áíçííæíúì.

1.3 Íáçíð íáèàñòè

Á ðááíòàò, ìòííñýùèòñý è íáèàñòè Polarization Detection è áèèçèíé è íáé Stance Detection, áàòíðù èñííèüçòðò ðàçèè÷íúá ìíñòáííèèè ðáðááííé çàää÷è: ìíñòáííèó ñ ó÷èòäèáí [13][6], ìíñòáííèó ÷áñòè÷ííá íáó÷áíèý [8] è ðáèääáííóð ááííé ðááíòá

īīñòàííáéó ááç ó÷èòáëÿ īīñēááíþþ ðàññīìòðèì īīäðíáíáá. Ðàçìáòèà àúáíðíē ðá-
ñòðñīçàòðàòííá äáéñòàèà. Áíēáá òíáí, áíēüøèíñòáí ñóúáñòáópùèò ðáøáíéé çàää÷è
Polarization Detection íà īñííáá ðàçìá÷áíúò íááíðíá äíēóíáíòíá Twitter èñīēüçò-
þò éííòðíēèðóáíúá īīäðíáú, éíòíðúá, èàé áúèí īīēàçáí, íá īīäàþòñÿ íáíáúáíèþ
è ñòðàääàþò ìò īíðáíè÷áííé äíñòóíííñòè é äáíáðáèüííé ñíâíēóíííñòè è é íáääæíúì
íááíðáì ááííúò, īīäòááðæààþùèò äíñòíááðííñòù īīēó÷áííúò ðáçòéüòàòíá [4].

Íñííáííé èääáé ðááíò, á éíòíðúò īðèíáíÿèíñú íáó÷áíéá ááç ó÷èòáëÿ, īíæíí ñ÷è-
òàòù ðáøáíéá çàää÷è īīēñèà èàòáíòííáí ááèòíðííáí īðíñòðáíñòáá ìàèíé ðàçìáðíí-
ñòè, á éíòíðíí áúáíáíáá īðíēçáíáèòù èèàñòáðèçàòèþ. Á ðááíòá [10] ñòðíēèíñú íáúáá
ááèòíðííá íðíñòðáíñòáí īīēòè÷áñèèò áçäëÿáíá īīēüçíáàòáèáé ñáòè Twitter è ááè-
òíð īīēòè÷áñèèò ñòíðíí īīñðááñòáí èñīēüçíááíèÿ éííèðíááíèÿ īīēüçíáàòáèÿìè
ñòíðíííèò íóáèèèàòèé ñ óèàçáíéáì ññúèèè íá íèò èàé ìòèèèèíá è ñīīñòàáèáíèÿ
èò íóáèèèàòèé ñ íóáèèèàòèÿìè īīēòè÷áñèèò ñòíðíí. Á ðááíòá [16] īðèíáíÿàòñÿ çà-
ñèóæèáàþùèé áíèíáíèÿ īīäðíá ñ íáó÷áíéáì īðááñòááèáíéé áçäëÿáíá ìàòíáíí ááç
ó÷èòáëÿ (unsupervised belief representation learning) è ÿòí īīñòðíáíéá èàòáíòííáí
īðíñòðáíñòáá ñáìè èáíóþò ááòíðú), ñ īīííúþ ááðèàòèíííáí ááòíèíáèðíáúèèèà
íáá ááòáíèüíúì íáíðèáíòèðíááíúì áðàòíí īīēüçíáàòáèáé è éííòáíòá ñáòè Twitter,
áðàò īīèñúááàòñÿ ìàòðèòáé ñíáæííñòè, ááñà éíòíðíé ñòðíèèèññú áéí÷íúì íáðàçíí. Á
ñðááíáíéè ñí íííáèìè àèáíðèòíáìè, íáèíòíðúá èç éíòíðúò ñíáòèòè÷áíú äëÿ ááííúò
Twitter, ÿòà īíááèü īíèàçúááàò íáèèò÷èéá ðáçòéüòàòù ñðááè īíäðíáíá íáó÷áíèÿ ááç
ó÷èòáëÿ.

Áñá ÿòè ðááíòù òàè èèè éíá÷á íáúááèíÿáò íáèè÷èá íá ÿòáíá íáó÷áíèÿ íáèíòíðíáí
ðàçìá÷áííáí ááèèèàòèíííáí íááíðá äíēóíáíòíá ñ àíðèíðíí çàääáíúì èç ááííúò
÷èñèíí òáíòðíá īīēÿðèçàòèè, éíòíðíá áúñòóíááò á èà÷áñòáá ááèíñòááííí ááðííáí
áèíáðíáðáíàòðá èèàñòáðèçàòèè èèè ìò éíòíðíáí óóíèèíáèüíí çàáèñÿò ÿáíúá áè-
íáðíáðáíàòðú àèáíðèòíáí. Íññēááíèá ááá ðàññīìòðáíúá ðááíòù īíæíí īīñòàáèòù
á òáðíèíáò íáó÷áíèÿ ááç ó÷èòáëÿ ñ īðèáèèááèðíááííé éíòíðíáèèáé. Çà÷áñòóþ áá-
òíðú ñòàòáé äëÿ ááèèèàòèè íááèáé èñīēüçòþò ìòèðúòúá ááííúá èç ñáòè Twitter, á
íáíé èç ðááíò èñīēüçòþò ÿáííá íááèèèðíáíéá ìòííáíèÿ é īīēòè÷áñèèò ñòíðííá
á òðáòííèÿðííé īīēòè÷áñèíé ñèñòáíá.

Çàää÷àPolarization Detection äëÿ óõññêíâí ÿçûêâ â ïñòàííâêâ ááç ó÷èòáëÿ ñ ïð-
 âèèââèðíâáííé èíîíðíàòèèé ðâøàèèñû â ðááíòá [26], â ÿòíé ñòàòóâ ïðèááâèè è ïí-
 ñòðíáíèð ïðèçíàèíâíâí ïðíñòðáíñòââ ñ ïííùùð ïíââèè ìâøèâ ñèíâ íáâ òðèíèâòàìè
 ¼ñóáúâèò-ïðáâèèâò-íáúâèò, äëÿ ñèíòàèñè÷-âñèèò ðíèáé ñèíâ è èò ÷-àñòáé ðá÷è, òííàèù-
 ïñòÿìè è ñáìàíòè÷-âñèèèè ðíèÿìè ñèíâ ïí Òèèíðó [12] è ïíñèââòðóââí ïíñòðíáíèÿ
 ÿÿâèé èèâñòâðèçàòèè ïíððáñòâí òáìàòè÷-âñèé ïíââèè ARTM [25]. Äëÿ ïíèó÷-áí-
 íúò òáìàòè÷-âñèèèò áâèòíðíâ ðàçíàðííòè 2 (èííííáíòâ ñíòââòñòâóâò áâðíÿòííòè
 ïðèíâèèâèííòè èâèííó-òí ïíèðñó) ïðíèçâíèèñÿ ïâðáòíâ è æâñòèé èèâñòâðèçàòèè
 ïíððáñòâí ïíâðáòíðâ Argmax. Áâòíð ñíñòââèèè áââ ñíáñòâáííúâ òáìú ñ áâòíÿ
 ïíèÿðèçíâáííúè ïáíèÿìè â èâæâí è ïíèó÷èèè F1-íáðó â ðáèíá 0.85 äëÿ íáâèò
 òáì. Áñèè èñíèüçíâòú ïòàòèð äëÿ áèíàðííé ìàððèòú ïðèáíè [11], òíðíóéò ÿòíé
 ìàððèèè ïíæíí áúíèñàòú òàè:

$$F1 = \frac{P \cdot R}{\frac{P + R}{2}}$$

ðâøáíèâ çàää÷èPolarization Detection ìàðíâàìè íáó÷-áíèÿ ááç ó÷èòáëÿ ñ ïðèâèèâ-
 âèðíâáííé èíîíðíàòèèé äëÿ èçíà÷-àèúíí íáòèèñèðíâáííâí ÷-èñèâ ïíáíèè â òáìàò
 íâ ðâøàèèñû ìè äëÿ áíèèèèñèíâíâí ÿçûêâ, ìè äëÿ óõññêíâíâí.

1.4 Öâü ðááíòú

Áâáííé ðááíòâ òðááóâòñÿ ñðâáíèòú ðàçíúâ ïíâòíâú è ïñòðíáíèð èèâñòâðèçàòèè
 áúáíðíè ñ íáòèèñèðíâáííú ÷-èñèí ïíèÿðèçíâáííúò ïíáíèè, äëÿ èíîíðíè èèâñòâðú
 áúñòóíâèè áú â èà÷-àñòââ ïíèðñíâ ïíèÿðèçàòèè. Òðááóâòñÿ òàèæâ ñðâáíèòú ìàðíâú
 íâ ñíñíáííòú áúââèÿòú øóí â òáìàò.

Áòíðíé òâèùð èññèââíâáíèÿ ÿâèÿâòñÿ ïñòðíáíèâ ìàðíââ ïòáíèâáíèÿ ïíèó÷-áííúò
 ïíââèè èèâñòâðèçàòèè ìà ñíñíáííòú ñòðíèòú ðàçáèáíèâ òáì ìà ïííæâñòââ, ïòíâ-
 ðâæâðóèâ ïíèðñú ïíèÿðèçàòèè, è ìà ñíñíáííòú áúÿâèÿòú øóííáíè èèâñòâð.

1.5 Ìñòàííâêâ çàää÷è ìàøèííâí íáó÷-áíèÿ

Èíââòñÿ ñíâíèóíííòú òáì $X = f X_g^{\%}$, èâæââÿ èç èíòíðúò ïíèñáíâ ìàððèòáé
 íáúâèòíâ-ïðèçíàèíâ, è èò ðàçíàòèâ $Y = f Y_g^{\%}$. Ñòââÿòñÿ%íáçàâèñèíúò çàää÷è íáó-

÷áíèÿ áâç ó÷èòáëÿ ñ ïðèàèèääãèðíâáííé èíðíðíàòèèáé: íáíáðíäèì ïíñòðíèòù ïí-
 äáëü, ïðèíèìàðóð íà äðíäâ áúáíðéó $X_8 = f_{\mathbb{G}_j \mathbb{G}_j} \mathbb{2} \mathbb{R} \leftarrow 9 = \overline{1-j} \cdot g \mathbb{2} X$, à íà
 áúðíäâ ïðíèçáíÿÿòð áâ èèàñòáðèçàòèð, ò.á. ïäáëü äíèæíà áúääâàòù ïíæáñòâí
 $Y_8 = f_{\mathbb{H}_j \mathbb{H}_j} \mathbb{1} \leftarrow \mathbb{8} \leftarrow \mathbb{9} = \overline{1-j} \cdot g$

Äëÿ èçíäðáíèÿ èà÷áñòââ ïíèó÷áííé ïíääèè èñííèüçóáòñÿ ïíæáñòâí Y_8 . Òàè
 èàè ñðàáíèèàòüñÿ áóáò ðàçèè÷íúà àèáíðèòíü èèàñòáðèçàòèè, òðááóáòñÿ áúáðàòù
 èàèèà-òí áíáðíèà ïàððèèè èà÷áñòââ ïí ïòííðáíèð è ìè. Èì èáñòóíàðò BCubed-
 òí÷ííòù (%), BCubed-ïíèííòà (') è BCubed-F-íàðà () [1].

Íóñòù èíèè÷áñòâí äíèóíáíòíâ íáèíòíðíé ðàçíá÷áííé òàíü - = f_{Gg} ðàáíà # ,
 èàæáííò íáúáèòó \mathbb{G} ïíñòàèèáíà ïàðèà èèáññà $\mathbb{H}_j = \overline{1} = \mathbb{e}$ ïàðèà èèàñòáðà: $g = \overline{1} < \mathbb{1}$ ïí-
 æáñòâí èèáññíâ íáíçíà÷èì . = f. $g \mathbb{1} \mathbb{8} = \overline{1} \cdot g$, ïíæáñòâí èèàñòáðíâ = f $g \mathbb{1} \mathbb{8} = \overline{1} < g$.
 Áääääí òóíèòèð $2^1 \mathbb{G} \leftarrow \mathbb{G} = 1 \gg \mathbb{H}_j = \mathbb{H}_j \mathbb{1} \gg \mathbb{8} = : g \mathbb{1}$ Òíäââ ïàððèèè èà÷áñòââ áúðàæàðòñÿ
 ñèääòðùèè òíðíóèàè:

$$\begin{aligned} \% &= \frac{1}{\#} \frac{\tilde{O}}{\mathbb{G}_2} \cdot \frac{1}{j} \frac{\tilde{O}}{:\mathbb{8}} \mathbb{2}^1 \mathbb{G} \leftarrow \mathbb{G} \\ ' &= \frac{1}{\#} \frac{\tilde{O}}{\mathbb{G}_2} \cdot \frac{1}{j \cdot \mathbb{H}_j} \frac{\tilde{O}}{\mathbb{G}_2 \cdot \mathbb{H}_j} \mathbb{2}^1 \mathbb{G} \leftarrow \mathbb{G} \\ &= 2 \frac{\%}{\%_0} \frac{'}{'} \end{aligned}$$

Íàððèèà BCubed-F-íàðà óáíäèèàðíðÿàò ÷áòùðáì ñáíèñòââì, ÿáðèñòè÷áñèè òàðàè-
 òáðèçòðùèì %ðíðíðòðç èèàñòáðèçàòèð ïíðááñòâí ïòííèòáèüííèòè ðàçóèüòàòíâ.

Íòääèüíí ïðíèííáíòèðòáì ñáíèñòââ %íáðèà ñ ïóííðííç è %ðàçíáð èèàñòáðà ïðí-
 òèà èò èíèè÷áñòââç: ïàððèèà èà÷áñòââ äíèæíà áúòù ÷óáñòáèòáèüíà è äðóííèðíâáíèð
 ðàçíáúáííü %íóííðííüç òí÷áè áúáíðèè â íäèí èèàñòáð è è ìèíèèçàòèè ÷èñèà
 èèàñòáðíâ â ñèó÷áâ, áñèè ÿòí íá äðáàèò íáúáíò íæèääíèð ïò ðàçóèüòàòíâ èèàñòáðè-
 çàòèè, ÷òí ñòíæèà íáúáèòù ïíáúàðòñÿ â íáíí ïíæáñòâí, à íáííòíæèà â ðàçíúâ.
 Íàððèèò BCubed-F-íàðà áóáâ ñ÷èòàòù ïíííáííé ïàððèèé èç áúáðáííü.

Íàððèèà BCubed-òí÷ííòù èíèääèèòñÿ â ïðííáæóðèèà ïò $\frac{\leq}{\#}$ äí 1, BCubed-ïíèííòà
 ïò $\frac{\equiv}{\#}$ äí 1, BCubed-F-íàðà ïò $\frac{2}{\#} \frac{\equiv}{\equiv}$ äí 1.

Def. 2: Nafēnōāā iāòðèèè èà-āñōāā èèāñòāðèçàöèè^{1 0}

Íáíçíà÷èì . 1 çà ¼øóíîáíéç èèāññ. Äëÿ ñðāáíáíèÿ ïîāāèāé ïā ñîñîíáííñòù áûāā-
 èÿòù øóì áóāāì èñîñèüçíāàòù àèáíèçíáíáíúā èñîñíáíúā ïāòðèèè. Áāāāāì #>8B4%
 $\frac{1}{j \cdot 1j} G_2 \cdot 1 \cdot 8 G_2 \cdot 8$ 2'G- G', #>8B4' = $\frac{1}{j \cdot 1j^2} G_2 \cdot 1 G_2 \cdot 1$ 2'G- G'. Íōáíéíé ñîñîñíáííñòè àè-
 āíðèòìā èèāñòāðèçàöèè áûāāèÿòù øóì ïāçíāāì BCubedNoise-F-íāðā (#>8B4):

$$\#>8B4 = 2 \frac{\#>8B4\% \#>8B4'}{\#>8B4\% \#>8B4'}$$

Ìāòðèèā ïðèíèìāàò òāì áíèüøèā çíà÷áíèÿ, ÷āì èó÷øā óāāèíñü ïîññòèòù øó-
 ïāúā íáúāèòù ā ïāèí èèāñòāð, ïðè ÿòì ïā ïòíāñÿ è íáíó íāøóíîāúā íáúāèòù.
 Ìāèñèìāèüííā çíà÷áíèā 1 ñîññāñòñòāóāò èāāāèüííé èèāñòāðèçàöèè äëÿ øóííāúō äí-
 èóíáíòíā.

Èà-āñòāí ïîāāèè, èāè ā ñíüñèā ñîññíáííñòè ñððíèòù èèāñòāðèçàöèè òāìù, òāè è
 ñîññíáííñòè áûÿāèÿòù ā ïèð øóì, áóāāì ñ÷èòāòù èāè ñðāáíāā èà-āñòāí ā èāæáíé òāìā.
 Ìðè ÿòì äîññèíèòāèüííí áóāāì ðāññíāòðèèāòù ïòāāèüíí ñðāáíāā èà-āñòāí ā ñðāçāò
 ÷èñèā ííáíèè ā òāìāò:

^ 1 ííáíèā;

^ 3 ííáíèÿ;

^ 2 ííáíèÿ;

^ 4 ííáíèÿ è áíèüøā;

2 Ðàøáíèà

Á ýòíé ñàèöèè ðàçááðàì ñííííá ìíñòðíáíèý ðàçíà÷áííúð àúáíðíé è ñííííá ðàø-
íèý ìíñòààèèáííé çààà÷è ìàøèííáì íáó÷áíèý.

2.1 Êèàñòàðèçàöèý

Ìðèíöèèèèèíá ìèèè÷èà ìíñòààèèáííé çààà÷è ìò òáð, ÷òì ðàøàèèñù à ñòíæèð ðà-
áíòàð, ýàèýàòñý íàèè÷èà íáðèèñèðíááííáì ÷èñèà ìíýðèçíááííúð ííáíéè à òáìàð.

Áóáàì ñ÷èòàòü, ÷òì èàæäüé ýèáíáíò èñòíáííé àúáíðèè ìíñòààèèáííé íáèíòíðüì
áàùáñòàáííüì áàèòíðì ðàçíàðííòè . Á ýòíé ðàáíòà ðàññííòðè ìàçíáúà àèáíðèò-
íü áàèòíðíé èèàñòàðèçàöèè.

2.1.1 Áèáíðèòíü, ìíííááííüà íà áíàèèçà íèíòííòíüð òàðàèòàðèñòèè àúáíðèè

Ìáðáüé ðàññííàðèèááííüé àèáíðèòíü $\frac{3}{4}$ Ìíííááííüà íà íèíòííòè ìðííòðáí-
ñòááííü èèàñòàðèçàöèý äèý ìðèèíæáííé ñ øíàìè ζ (DBSCAN) [9]. Èíáàò
òðè èèð÷ááüð àèíáðíàðáìáòðà:

^ eps ðàññòíýíèà, ìðáááèýðáá ñííááñòáí òí÷áè (ááá òí÷èè íáýýàèýðòñý ñííáá-
íèè, áñèè ðàññòíýíèà ìæäó íèè ìáíüøá eps);

^ min_samples: ìèíèèèèíá èíèè÷áñòáí òí÷áè äèý ìðáááèèèè èèàñòàðà;

^ òóíèöèý ðàññòíýíèý/ìàðèèè : èðááý òóíèöèý, óáíáèáòáíðýðáý àèñèíáì
òíæááñòáà, ñèìàðèè è òðáóáííèèèè;

Ìà èð ìíííáá òí÷èà ìðííòðáííòáà ìíæàò áúòü ìðáááèèà èàè:

^ ýäðíááý òí÷èà : áíèðóá íáá á ðàèèóñáeps ìí çàááííé ìàðèèè ìàðíäýòñý
min_samples òí÷áè, ò÷èòüááý ñàìó ýäðíáóð òí÷èò;

^ áðáíè÷íü òí÷èà : á ðàèèóñáeps ìò íáá áñòü òíòý áú íáíà ýäðíááý òí÷èà, ìí á
òí æá ðàèèóñá íàò ðàññííèèèè ìáíüøá min_samples òí÷áè ñ ò÷áòí ñáííé
áðáíè÷íè òí÷èè;

^ áúáðíí: áí áñáð ìíñòàèèíüð ñèó÷áýð;

Ōiöy äëäíðèèì ñääëää ñðääíáái çíà÷áíèy øèðíèí èñííèüçóáòñy âí íííãèõ ïðèèí-
æáíèyõ, ñòðííáí áíèàçàòáèüñòái ñííàèíñòè àèäíðèèì, èñííèüçóðùääí yäðí íáúääí
âèää â ïðíñòðáíñòääò âñííèíè ðàçíàðííñòè, ïòñòòñòâóâò.

Íááíð ðàññíàòðèääáíüõ äèíáðíàðàìàòðíá:

- ^ ðàçíàð íèíá h: » 20-0¼íí èííàðèèòè÷áñíèé ñàòèá ïí ïíííáàíèð 4;
- ^ íííæáñòái öáíòðíèäíá : áñá òí÷èè âúáíðèè;
- ^ óóíèöèy yäðà : ¼íèííèíáç yäðí, çàääääáííá òíðíóèíè $1\mathfrak{G} = 1 \gg \mathfrak{J}\mathfrak{J}_2$ ¼

2.1.2 Áèäíðèèìü, yáíí ïàðàìàòðèçóáíüá èíèè÷áñòái ìèèñòòðíá

Íáðáúé ðàññíàòðèääáíüé àèäíðèèì ñíáñü ííðíàèüíüõ ðàññíðääáèáíèé
Ñòðíèòñy áííðíèñèìàöèy ðàññíðääáèáíèy èñííáííè áúáíðèè - ïíííñòè # íáèíòí-
ðúi ïíáèèüíü ðàññíðääáèáíè, ñíáñüð íííííáðíüõ ííðíàèüíüõ ðàññíðääáèáíèé
ñ K èííííáíòàìè è àíðèíðíüè ááñáìè c: , çàääääáííè ñèááóðùèè òíðíóèàìè:

$$?^{1-0} = \begin{matrix} \mathfrak{O} & \mathfrak{O} \\ c: N^1 \mathfrak{G}_j \backslash : - :^0 - & \mathfrak{O} \\ \mathfrak{g}=1 : =1 & : =1 \end{matrix} c: = 1 - \mathfrak{g} \quad 0$$

Áääääí íáíçíà÷áíèy c = f c_{1-...-g}, ` = f`_{1-...-g}, = f_{1-...-g}. Íðääíí-
èääääòñy, ÷òí èàæáúé íáúáèðGáúáíðèè ïðèíàèèæèò íáííè è òíèüèí íáííè èíííí-
íáíòá ñíáñè. Áííàèòñy áðóííà ñèðùòüõ ïáðáíáííüõ / = f|_{1-...-g}, ïíèñüâàðùèò
ïðèíàèèæííñòü èàæáíí íáúáèòà íáèíòíðíè èííííáíòà, |_{8 2 1-} - ïíáð èíííí-
íáíòü ñíáñè, èíòíðíè ïðèíàèèæèò íáúáèòà G. Áíðèíðíá ðàññíðääáèáíèá áàèòíðà
|<sub>8 2 / : ?^1|_g c^0 = Cat^1_g c^0 = \begin{matrix} 1 \\ c: \cdot \\ \cdot \\ \cdot \\ \cdot \end{matrix} . ïàðàìàòðü ïíáèèè f c: - \cdot - : g_{=1} íáñòðàèèàðòñy
ïðè ïííüè EM-àèäíðèèì [7] Áèíáðíàðàìàòðíá àèäíðèèì áññòóíàðò:</sub>

- ^ ÷èñíè èííííáíò ñíáñè;
- ^ f c-`-g íà÷àèüíáy èíèèèèèçàöèy ïàðàìàòðíá ïíáèèè;

Áóääí èñííèüçíáàòü ñèó÷áèíòð èíèèèèèèçàöèð ïàðàìàòðíá ïíáèèè, èó÷øäy ñðá-
äè íáó÷áííüõ EM-àèäíðèèì ïíáèèè áúáèðääòñy ïí áàèè÷èíá áàðèàòèíííè íèæ-
íáè ïòáíèè.

Íááíð ðàññíàòðèääáíüõ äèíáðíàðàìàòðíá:

$$\hat{\mu} = \arg \min_{1 \leq j \leq K} \sum_{i=1}^n \|x_i - \mu_j\|^2$$

$$\hat{\mu} = \arg \min_{1 \leq j \leq K} \sum_{i=1}^n \|x_i - \mu_j\|^2 \quad (40 \text{ \textit{K-means}});$$

Àðìðíé ðàññìàòðèääàíúé àèáíðèòì ìàòíà K-ñðääíèõ (K-means) [22]. Ìðää-ñòàâëÿàò ñíáíé ìðääáèüíúé ñèó÷-àé ìñòðìáíëÿ ñíàñè ìðìàèüíúò ðàññìàòðèääàíéé ìðè àçÿòèè ðàáíúìè àìðèìðíúò àáðìÿòìñòàé èììííáíò, òèèñèðìàáíèè àèàáííàèüííáí àèää ìàòðèò èíààðèàòèè : è ñòðàìèáíèè àá èììííáíò è ááñèííá÷íñòè. Àèìàðìà-ðàìàòðàìè àèáíðèòìà ÿàèÿðòñÿ:

$$\hat{\mu} = \arg \min_{1 \leq j \leq K} \sum_{i=1}^n \|x_i - \mu_j\|^2$$

$$\hat{\mu} = \arg \min_{1 \leq j \leq K} \sum_{i=1}^n \|x_i - \mu_j\|^2 \quad \text{èç èìòìðúò ñòàðòòòò èòàðàòèáíúà ìðì-òáññù ìñèñèà òáìòðìà èèàñòàðìà};$$

Èñìíèüçòÿ àáääáííòð àèÿ ñíàñáé ìòàòèò, ìñèáì çàääàòù èòàðàòèèííúé ìðìòáññ àèáíðèòìà ñèääòòèè òìðìóèàìè:

$$\mu_j = \arg \min_{\mu} \sum_{i \in C_j} \|x_i - \mu\|^2 \quad j = 1, \dots, K$$

Ìááíð ðàññìàòðèääàíúò àèìàðìàðàìàòðìà:

$$\hat{\mu} = \arg \min_{1 \leq j \leq K} \sum_{i=1}^n \|x_i - \mu_j\|^2$$

$$\hat{\mu} = \arg \min_{1 \leq j \leq K} \sum_{i=1}^n \|x_i - \mu_j\|^2 \quad \text{K-ñðääíèõ} \quad [2];$$

Ìñèääáíéé ðàññìàòðèääàíúé àèáíðèòì áàèáñííàñèàÿ ìðìàèüíàÿ ñíàñù [18]. Ìñìíáíé èääáé èðáíé áàèáñííàñèé ìñèèè ÿàèÿàòñÿ àáääáíéà àìðèìðííáí ðàññìàòðèää-èáíëÿ ìà ìðàìàòðù ìñèèè. Ááääáì ìáðçìà÷-áíèà = f 1-...-g, áää : = . 1. Ìáðàòìàÿ ìàòðèòà ñóúáñòàóàò â ñèò ìáâðìæääííñòè ìàòðèòù èíààðèàòèè ìð-ìàèüííáí ðàññìàòðèääàíéé. Áóááì èñìíèüçìààòù ìàðàìàòðù àìàñòì ìàðàìàòðìà , áääü ì ìáðàùì ìàðàìàòðàì àçàèìí-íáíçìà÷-ííèó÷-àòòñÿ àòìðúà.

Ìðääà òàì, èàè ñòðìèòù àáðìÿòìñòòòò ìñèèè áàèáñííàñèé ìðìàèüííé ñíàñè, ááääáì ìáñèííèè ðàññìàòðèääàíéé. Ááääáì ðàññìàòðèääàíéà ìàà ñèìàòðè÷-íúìè ìñè-æèòàèüíí ìðääáèáíúìè ìàòðèòàìè.

Íónòü 2 R³ =) < 0, , = ,) < 0, a 3 1. Òĩãã ðàñĩðãããããããã

$$?^1 j^-, 0 = 1, -a^{01} \det^{0 \frac{a-3-1}{2}} \exp^1 \frac{1}{2} \text{atr}, 10-$$

ããã 1, -a⁰ íðìèðíáí÷íàÿ éííñòàíòà ðàñĩðãããããããã, áóããí íàçúããòü ðàñĩðãããããããã-íèàì Óèðàððà.

Íónòü ` 2 R³, < 2 R³, V j 0. Òĩãã ðàñĩðãããããããã

$$?^1 -_0 = NW^1 -_j < -V - a^{-0} = N^1 -_j < -V^0 1^0 W^1 -_j a^{-, 0}$$

áóããí íàçúããòü Óèðàðð-íðìàèüíúì ðàñĩðãããããããã.

Íónòü c 2 R - 8 0- 8 c 8 = 1, U 2 R - 4 j 0. Òĩãã ðàñĩðãããããããã

$$?^1 c j U^0 = \text{Dir}^1 c j U^0 = \frac{1^1 U_8^{00} \ddot{O}}{1^1 U_8^0} c_8^{U_8 1}$$

íàçúããòü ðàñĩðãããããããã Æèðèõèã.

Ñ ó÷: àòìí áããããããã ãìðèíðíúõ ðàñĩðãããããããã ìà ìàðàìàððü íðìàèüíé ñìãñè
f c: - : - : g_{=1} áãðíÿòííñòíóð ìãããü ìæíí çàìèñàòü ñèããóðüèì íàðàçì:

$$?^1 -_ / - c - - 0 = \frac{\ddot{O}}{8=1} ?^1 G_j |_{8-} - 0 ?^1 |_{8j} c^0 ?^1 - 0 ?^1 c^0$$

$$?^1 G_j |_{8-} - 0 = N^1 -_ |_{8-} |_{8}^{10}$$

$$?^1 |_{8j} c^0 = \text{Cat}^1 |_{8j} c^0$$

$$?^1 -_ |_{8-} |_{8}^0 = NW^1 -_ |_{8-} |_{8j} < 0 - V - a^{-, 0}$$

Ðàñĩàòðèããòñÿ 2 òèìà ááéãñíãñèéé íðìàèüíúì ñìãñé, ìò ìñòàííãèè çããèñèò çãããíèãã ãìðèíðíá ðàñĩðãããããããã ìà c:

$$\hat{\ } \text{éííá÷íàÿ: } ?^1 c j U_0^0 = \text{Dir}^1 c j f U_1 - \dots - U_4^0,$$

$$\hat{\ } \text{áãñèíá÷íàÿ: } ?^1 c j U_0^0 = \text{Dir}^1 c j f U - \dots - U g^0, \text{ããã } ! 1 ;$$

Äëÿ éííá÷íé ñìãñè ìðèìáÿòñÿ ñòàìà Æéããñà äëÿ ãáíáðàòèè ñéó÷: áéííáí áãèòíðà èç ðàñĩðãããããããã Æèðèõèã [14], äëÿ áãñèíá÷íé ñìãñè ñèíèèàíèðíááíáÿ ñòàìà Æéããñà [17].

Ἄ αἰίίε δααίòδᾶ αóαᾶò èñîîëüçîáàòüñý ðᾶᾶèèçàòèý ᾶᾶñéíᾶ-ííé ñìᾶñè ííðìᾶèüíüð
 ðᾶñîðᾶᾶᾶéᾶíéè ÷ᾶðᾶç ᾶîîðíēñèìᾶòèþ ñòᾶíü ᾶᾶíᾶðᾶòèè ñ îîîüüþ îðíòᾶññᾶ ¼éííèè
 ìᾶèèèèè (Stick-breaking) [3]. Ἐèþ-ᾶᾶüì ᾶèìᾶðìᾶðᾶìᾶòðîì ᾶèᾶîðèòìᾶ ᾶüñòóìᾶᾶò ìᾶðᾶ-
 ìᾶòð U₀ ðᾶñîðᾶᾶᾶéᾶíéý ?¹c j U₀⁰, ñîñòîýýùè èç K ìᾶèìᾶéíᾶüò ìîéíᾶèòᾶèüíüð ᾶᾶèè-
 ÷éí. Çᾶᾶᾶíèᾶ ìèçèíᾶí çìᾶ-ᾶíéý ìðèᾶíᾶèò è éííòᾶíòðèðìᾶᾶèþ ᾶᾶðíýòíñòóíé ìᾶññü
 ᾶ ìᾶᾶíèüðé ÷èñᾶ éíîîíᾶíò ñìᾶñè è ðᾶçðᾶᾶèᾶᾶèþ ìᾶᾶèè ÷ᾶðᾶç ìðèᾶèèᾶᾶíᾶ è ÷ᾶ-
 ñòè ᾶᾶñᾶ éíîîíᾶíò è 0. Óᾶᾶèè-ᾶíᾶ çìᾶ-ᾶíéý ᾶèìᾶðìᾶðᾶìᾶòðᾶ ìçᾶíᾶýᾶò ìñòᾶᾶèòü
 ᾶ ìᾶᾶèè ᾶíèüøᾶᾶ ÷èñᾶí èíîîíᾶíò.

Äèý èíèèèᾶèèèèèèè ìᾶðᾶìᾶòðìᾶ ᾶóᾶᾶì èñîîëüçîáàòü ñèᾶᾶóþùèᾶ ᾶᾶèè-èíü:

$\hat{c} <_0$: ñðᾶᾶíéè ᾶᾶèòð óᾶíòðìᾶ èèᾶñòᾶðìᾶ, ìîéó-ᾶíüò ñ îîîüüþ ᾶèᾶîðèòìᾶ Ἐ-
 ñðᾶᾶíéò;

\hat{V}_0 : 1;

\hat{a}_0 : ;

\hat{c} , c_0 : ìᾶòðèòᾶ éíᾶᾶèèèèèèèèèè ᾶèý ìᾶòðèòü ìᾶúᾶèòìᾶ-ìðèçìᾶéíᾶ ᾶúᾶíðèè;

Íᾶᾶíð ðᾶññìᾶòðèᾶᾶᾶüò ᾶèìᾶðìᾶðᾶìᾶòðìᾶ:

\hat{c} : »1-minf 25-j- jg¼

\hat{c} èíîîíᾶíòü ᾶᾶèòðᾶ U₀: » 6- 1¼ñ éíᾶᾶðèòèè-ᾶñᾶíé ñᾶòèᾶ ìñ ìñíᾶᾶèþ 4;

2.1.3 Ἄúᾶíᾶú

Ἀèᾶîðèòì ìñòðìᾶíéý ᾶᾶéᾶñìᾶñᾶíé éíðìᾶèüíé ñìᾶñè δᾶᾶíòᾶᾶò ᾶíèüøᾶ ìñòᾶèüíüð,
 ìí ìñçᾶíᾶýᾶò ìᾶᾶèèðìᾶᾶòü ìᾶíüøᾶᾶ èçìᾶ-ᾶèüíí çᾶᾶᾶííᾶí ÷èñᾶí èíîîíᾶíò ñìᾶñè çᾶ
 ñ-ᾶò çᾶíóéᾶíéý ᾶᾶðíýòíñòᾶé ñîòᾶᾶòñòᾶóþùèò ᾶᾶñᾶ éíîîíᾶíò. Íí ᾶñᾶ ᾶúᾶ òðᾶᾶóᾶò
 òíᾶíᾶí ìᾶᾶíðᾶ ìᾶðᾶìᾶòðᾶ ðᾶñîðᾶᾶᾶéᾶíéý Ἀèðèðèᾶ.

Ἐᾶᾶᾶüé ᾶèᾶîðèòì èèᾶñòᾶðèçᾶòèèè èðᾶᾶóᾶò ìᾶñòðìᾶéèè ìᾶéíòðìᾶí éíèè-ᾶñòᾶᾶ ᾶè-
 ìᾶðìᾶðᾶìᾶòðìᾶ ìᾶíᾶóᾶíᾶèí èñîîëüçîáàòü èᾶéíé-òí èðèòᾶðèè ìᾶíðᾶ ìᾶᾶèè ᾶèý èò
 ìᾶᾶíðᾶ. Ἄ èᾶ-ᾶñòᾶᾶ òᾶéíᾶíᾶí ᾶóᾶᾶì èñîîëüçîáàòü éíýòèèèèèèèèè ñèèóýòᾶ [21]. Ðᾶñ-
 ñîòðèè ᾶúᾶíðèò - ìüíñòè #, ñîñòîýýóþ èç ìᾶúᾶèòìᾶ G2 R . Ἐᾶᾶᾶííó ìᾶúᾶèòó
 ìðèèèñᾶíᾶ ìᾶòèᾶ èèᾶñòᾶðᾶ H2 1- , ìóñòü = f 1-•••- g ìíᾶᾶñòᾶí èèᾶñòᾶðìᾶ,

íáðàçòðùèò ðàçáèàíèà áùáíðèè. Çàòèèñèðóàì íáèìòíðòð íàòðèèó $d^1 -^0$. Äëÿ èàæ-
 äíáí íáúáèòà $G_2 -$, äëÿ èìòíðíáí $j_{H_j} j - 1$, ïðääääèè:

$$\hat{0}_8 = \frac{1}{j_{H_j} - 1} \int_{G_2} d^1 G_2 - G_2 \quad \text{ñðääíáá ðàññòíÿíèà ò èñòíáííé òí÷èè áí áñáò}$$

òí÷áé òíáí æá èèàñòáðà;

$$\hat{1}_8 = \min_{H_k H_b} \frac{1}{j_{H_j} - 1} \int_{G_2} d^1 G_2 - G_2 \quad \text{íàèíáíóðáá ñðääíáá ðàññòíÿíèà ò èñòíáííé òí÷èè áí}$$

áñáò òí÷áé äðóáíáí èèàñòáðà;

Êíÿòèèèèáíò ñèèóÿà íáúáèòà G_2 áùáíðèè ïðääääèè èèè

$$B_8 = 1 \gg j_{H_j} j - 1 \frac{1}{\max f 0_8 - 1} \frac{0_8}{1g}$$

ÿòà ááèè÷èíà íàðíàèòñÿ á òðáçèá » $1 - 1/4$ è ïíèàçúááàò, íàñèíèùèí òíðíòí óää-
 èíñü èèàñòáðèçíáàòü èíéðáòíúé íáúáèò áùáíðèè. Áíèää èà-áñòááíáÿ èèàñòáðèçà-
 öèÿ èíáàò áíèüøáá çíà÷áíèà èíÿòèèèèáíòà ñèèóÿà, ïðè ÿòí òéääíá çíà÷áíèà áù-
 ñòóíáàò èíàèèàòíðíí íàñèíáíèÿ èèàñòáðíá, à òðèèòáèèúíá çíà÷áíèÿ èíàèèàòí-
 ðíí òíáí, ÷òí áíèüøèíðáí íáúáèòíá áùáíðèè èìáðò íáááðíòð èèàñòáðèçàòèð, ò.á.
 äëÿ áíèüøèíðáá íáúáèòíá èàèíé-òí äðóáíé èèàñòáð èç èìáðùèòñÿ áúè áú áùáí-
 ðíí èó÷ðá, ÷òí áúá ïæáò áíáíðèèòü ï ñèèøèí áíèüøíí èèè ñèèøèí íàèíí ÷èñèá
 èèàñòáðíá.

$$\text{Êíÿòèèèèáíò ñèèóÿà äëÿ èèàñòáðèçàòèè çàááàòñÿ èà} B_8 = \frac{1}{\#} \frac{\#}{8-1} B_8$$

Èñòíáííé èíÿòèèèèáíò ñèèóÿà íá ïðääääèèáí äëÿ èèàñòáðèçàòèè ñ íáíèì èèà-
 ñòáðíí áííðääääèèì ááí äëÿ ÿòíáí ñèó÷áÿ òéääúì çíà÷áíèàì. Õíááà èòíáíáúé
 èíÿòèèèèáíò ñèèóÿà ïðèíàò áèà:

$$B = 1 \gg j - 1 \frac{1}{\#} \frac{\#}{8-1} B_8$$

Òàèèì íáðàçíí, èó÷èè íááíðíí àèíáðíáðáíàòðíá àèáíðèèòíà èèàñòáðèçàòèè áó-
 ááí ñ÷èòáòü òàèíé, äëÿ èíòíðíáí èíÿòèèèèáíò ñèèóÿà íàèñèíàèáí.

Íàòðèèó äëÿ èíÿòèèèèáíòà ñèèóÿà áóááí áùáèðáòü òó æá, èíòíðòð èñíèèüçóàò
 àèáíðèèòí èèàñòáðèçàòèè, áñèè íàòðèèà áùñòóíáàò ááí àèíáðíáðáíàòðíí, èíà÷á èñ-
 ïèèüçóáì ááèèèáíáò íàòðèèó.

2.2 Áâèòíðííá ìðääñòàâèáíèá

Áâèòíðííá ìðääñòàâèáíèá áóääì ñòðìèòü ìáçàâèñèì äëÿ èàæäíé èç áúáíðíè. Á èà-áñòàá ààçíáíáí áâèòíðííáí ìðääñòàâèáíèáí áíèóíáíòíá ìðèíèàáòñÿ ñèääòpùèé ìááíð ìðèçíàèíá:

- ^ èíèè-áñòáí ìðèääááííü è ìà-àèüííé òíðíá ìðòíèíáè-áñèè ìáèèçàòíðíì pymorphy2 [15] èìáíááííü ñóúíñòáé, èìáðùèò ìíèíæèòáèüíáÿ, ìáèòðàèü-íáÿ è ìòðèòáòáèüíáÿ òííáèüííñè á áíèóíáíòá, áàç ó-áòà ìðÿäèà ñèíá á ìðá-áíáðááíòáííé èìáíááííé ñóúíñòè (ðàçíáðííñòü áàðüèðóáòñÿ ìò 1 áí 3);
- ^ ñóíàðííá çíà-áíèá òííáèüííñè äëÿ èàæäíé èìáíááííé ñóúíñòè, ìðèääáá-ííé è ìà-àèüííé òíðíá ìðòíèíáè-áñèè ìáèèçàòíðíì, èíòíðáÿ áñòðá-áòñÿ áí áñáò òáèñòáò áúáíðèè òíòÿ áú 2 ðàçà áàç ó-áòà ìðÿäèà ñèíá á èìáíááííé ñóúíñòè (ðàçíáðííñòü áàðüèðóáòñÿ ìò 9 áí 55);
- ^ èíèè-áñòáí èìáíááííü ñóúíñòáé ìðääáèáííé èàòáíðèè (¼èè-ííñòü, ¼ááí-áðáòè-áñèíá ìíèíæáíèá, ¼íðááíèçàòèÿ) (ðàçíáðííñòü áàðüèðóáòñÿ ìò 1 áí 3);
- ^ áðóííá ñíòèàèüíí-ááííáðáòè-áñèè ìðèçíàèíá èç èñòíáíáí ìááíðá áíèóíáíòíá (ðàçíáðííñòü 15);
- ^ ÿíáääáèíá ìðááíáðááíòáííáí ìáúáàèíáíèÿ çááíèíáèà è òáèà áíèóíáíòá, ìíèó-áííúé ñ ìííüð ìðááíáó-áííáí SBERT [23], ñ ìáðáíè-áíèá ìá èíèè-áñòáí áòíáíüò òíèáíá 384, á ñèó-áà áñèè áíèóíáíò èìáò áíèáá èíðíòèíá áíóòðáííáá äëÿ ìáèðíííé ñàòè ìðääñòàâèáíèá, áðíáííé áâèòíð áíçáííèíÿáòñÿ ¼íóèááúíè, òíèáíáè (ðàçíáðííñòü 1024);

Ìðèääááí ìòèääòèð èñíèèüçíááíèÿ òáèèò ìðèçíàèíá.

1. Íáíèá í ñíáúòèè, á èíòíðíì èèð-ááúíè ¼ó-áñòíèèàìèç ìíáòò áúñòóíáòü èè-ííñè èè ìðááíèçàòèè, ááèñòáòpùèá á ìðääáèáííì ìáñòá, òíðíèðóáòñÿ ìòí-ðáíèáí è ÿòèì ñáíúí èè-ííñòÿ / ìðááíèçàòèÿ / ìáñòá ìðíèñòáñòáèÿ, á ÿòí ìòííðáíèá ìàòáàòè-áñèè ñòíðíóèèðíááíí á áèáá òííáèüííñè, ìðèèñáíííé èìáíááííé ñóúíñòè ñíðáòñòáòpùáé èàòáíðèè.

2. Ííáíéá, ìðíááèääáííá á ÑÌÈ ïíáíñòíúìè ðáñóðñàìè, íàìðàáèáíí íà ïíðááá-
èáííúé ñðáç èþááé, ááá ñðáçàìè ïíáóò áúòù íáðàçíááíéá, áíñòàòíè èèè ïíè.

3. Ííááèü SBERT íáó÷àèàñü ðáøàòù íáñêíèüèí çàää÷ ïíèìàìèý áñòáñòááííáí
ýçúèà íáííáðáííí (multitask-learning), à èìáíí:

- ^ çàää÷à ààòíìàòè÷áñêíáí ïíðááèááíèý èíáè÷áñêíé ñáyçè íáæáó òáèñòàìè
(natural language inference) èìáàòñý íáèíòíðáy òáèñòíááy ïíñúèèà, ïíáá-
èèðóþùáy ñèòóáòèþ, íáèíòíðíá òáèñòíáíá ìðááííèíæáíéá í ñèòóáòèè, çà-
ää÷à ñíñòíèò á èèáññèòèèáòèè ìðááííèíæáíèý íà òðè èèáññà: ¼ìðáááàç,
¼èíæüç, ¼ííðáááèèòù íááíçííæííç;
- ^ çàää÷à áúááèáíèý èìáííááííúò ñóúííñòáé;
- ^ çàää÷à àíàèèçà òííáèüííñòè èìáííááííúò ñóúííñòáé;

Ýòè çàää÷è áèèçèè è ðáøááííè á ýòíè ðááíòá çàää÷á ñ òí÷èè çðáíèý áúááèää-
áííè ïíòèääòèè èñííèüçíááíèý ìðèçíáèíá, ïíýòíó áñòù ïíáíá ïíèáááòù, ÷òí
ïíðáááèýáííá ïíááèüþ ááèòíðíá ìðíñòðáíñòáí ýíáááèèíáíá ïíçáíèèò ïíñòðí-
èòù èèáñòáòù, íòíáðáæàþùèà èìáííí ïíèýðèçíááííúá ííáíèý.

Ìðíèçááááí ñòáíáàððèçàòèþ ìàòðèòù íáúáèòíá-ìðèçíáèíá: áèý èáæáííá ìðèçíá-
èà íáçááèñèíí áú÷òáí ááí ñðááíáá àðèòíàòè÷áñêíá è ïíááèè ïíñèá ýòíáí íà ñòáí-
áàðòííá ìèèííáíèá.

2.3 Ñíèæáíèà ðàçíáðííñòè

Èòíáíááy ðàçíáðííñòù ááèòíðá áíñòàòí÷íí ááèèèà, ïíýòíó íáíáðíáèí ìðèááá-
íòù è ìàòíáà ñíèæáíèý ðàçíáðííñòè. Áíííèüçóáíñý ìàòíáí áèááíúò èííííáíò
(PCA) áèý èèíáííáí ñíèæáíèý ðàçíáðííñòè èñòíáííáí ááèòíðíáí ìðíñòðáíñòáá [19].
Ñóòù ìàòíáà ñíñòíèò á ïíèñèá òáèíáí ìðíñòðáíñòáá íáíúøáé ðàçíáðííñòè, èíòíðíá
ìèìèìèçèðóáò ïíðáðþ èíòíðíàòèè í áúáíðèá ñ òí÷èè çðáíèý àèñíáðñèè. æèñí èíí-
ííáíò áóááí ïíááèðáòù òáèíá, ÷òíáú ïíèò÷áííúá ìðèçíáèè ñíððáíèèè 99% èñòíáííè
àèñíáðñèè. Íóñòù $_1-\dots-_m$ ïíæáñòáí ñíáñòááííúò çíá÷áíèè ýííèðè÷áñêíé èíáà-
ðèàòèíííè ìàòðèòù íáúáèòíá-ìðèçíáèíá. Ííè áñá íáíòðèòáòáèüíúá á ñèéò íáíòðè-

$\min_{g=1, \dots, g} \left| \frac{g-1}{g-8} \right|$

Δαῖτα ἰσοπέδη ἀεὶ ἔχει PCA (kernel PCA, kPCA) ἀεὶ ἰσοπέδη ἔχει
 ἰσοπέδη [24]. Ἄρα ἀεὶ ἰσοπέδη ἔχει ἀεὶ ἰσοπέδη ἔχει ἰσοπέδη
 ἀεὶ ἰσοπέδη: $\frac{G}{H} = \frac{G}{H}$. Ἄρα ἰσοπέδη ἔχει ἰσοπέδη ἀεὶ ἰσοπέδη
 ἰσοπέδη ἔχει ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη
 ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη

2.4 Δαῖτα ἀεὶ ἰσοπέδη

Δαῖτα ἀεὶ ἰσοπέδη ἔχει ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη,
 ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη
 ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη
 ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη
 ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη

Δεῖν. 3: Ἐἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη

Ἄρα ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη
 ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη ἀεὶ ἰσοπέδη

êîèè÷-âñòâî äîéóíáíòíâ â òáìàð. Ñðááíáâ êîèè÷-âñòâî øóííáúð äîéóíáíòíâ óáúââ-
 áò ñ ðñòòî ÷-èñèà ííáíéé â òáìàð, ïðè ýòíí äëý òáì ñ íáíèì ííáíéâì øóííáúð
 äîéóíáíòíâ áíëüøá, ÷-âì ïñýðèçíâáííúð.

3 Âú÷-èñèèòáëüíúé ýèñíáðèìáíò

Äëý ýèñíáðèìáíòíâ èññíëüçíââèèñü ðáàèèçàòèè àèáíðèòíâ èèàñòáðèçàòèè èç
 áèáèèòáèè scikit learn [20].

Â èà÷-âñòââ äññíèèòáëüííé èíóíðíàòèè í èà÷-âñòââ èèàñòáðèçàòèè áóáâì èññíëü-
 çíââòü ìàððèéó "8BB >D=C ðàçííòü ìæäó êîèè÷-âñòâî èèàññíâ è êîèè÷-âñòâî
 èèàñòáðíâ, ïñòðíáííúð àèáíðèòíí, óñðááíáíáý ïí áñáì òáìàì.

3.1 Ðàçóëüòáòü ýèñíáðèìáíò

Ïðèáááâì ðàçóëüòáòü â àèáâ ñáíáííúð òááèèò ïí áñáì òáìàì è ïí èàæáííó èç
 íáíçíá÷-áííúð â ïñòáííáèâ çàââ÷-è ñðáçíâ.

Ìíâáëü		Ìàððèèâ èà÷-âñòââ				
		%	'		#>8B4	"8BB >D=C
DBSCAN	PCA	0.565	0.919	0.674	0.428	1.69
	kPCA	0.584	0.869	0.668	0.428	1.36
Ñðááíáâî ñââèèâ	PCA	0.959	0.267	0.397	0.611	-9.86
	kPCA	0.897	0.371	0.479	0.670	-6.20
Ííðíàèüíáý ñíáñü	PCA	0.637	0.572	0.563	0.588	0.33
	kPCA	0.692	0.553	0.564	0.588	-0.38
K-ñðááíéð	PCA	0.679	0.774	0.681	0.646	0.36
	kPCA	0.764	0.510	0.569	0.600	-1.73
Áàéáñíâñèáý ïíðíàèüíáý ñíáñü	PCA	0.667	0.825	0.698	0.625	0.58
	kPCA	0.768	0.513	0.571	0.626	-1.87

Òááèèò 1: Ðàçóëüòáòü ïí áñáì òáìàì

Модель		Метрика качества				
		%	'		#>8B4	" 8BB >D=C
DBSCAN	PCA	0.898	0.929	0.896	0.668	-0.07
	kPCA	0.918	0.856	0.866	0.668	-0.27
Среднего сдвига	PCA	1.000	0.150	0.244	0.359	-11.13
	kPCA	0.992	0.221	0.347	0.416	-6.93
Нормальная смесь	PCA	0.933	0.506	0.641	0.417	-1.13
	kPCA	0.947	0.448	0.577	0.393	-2.07
К-средних	PCA	0.966	0.667	0.765	0.428	-1.40
	kPCA	0.984	0.429	0.569	0.435	-2.60
Байесовская нормальная смесь	PCA	0.966	0.748	0.825	0.412	-1.27
	kPCA	0.992	0.446	0.582	0.405	-2.67

2:

1

Модель		Метрика качества				
		%	'		#>8B4	" 8BB >D=C
DBSCAN	PCA	0.574	0.936	0.700	0.377	1.15
	kPCA	0.583	0.908	0.693	0.377	0.96
Среднего сдвига	PCA	0.932	0.292	0.417	0.716	-7.26
	kPCA	0.926	0.305	0.446	0.738	-6.44
Нормальная смесь	PCA	0.671	0.584	0.601	0.695	-0.33
	kPCA	0.733	0.547	0.596	0.728	-0.93
К-средних	PCA	0.733	0.769	0.724	0.739	-0.44
	kPCA	0.761	0.569	0.599	0.710	-1.52
Байесовская нормальная смесь	PCA	0.722	0.808	0.739	0.732	-0.26
	kPCA	0.766	0.565	0.595	0.655	-1.89

3:

2

Модель		Метрика качества				
		%	'		#>8B4	" 8BB >D=C
DBSCAN	PCA	0.493	0.892	0.627	0.445	1.76
	kPCA	0.497	0.865	0.621	0.445	1.69
Среднего сдвига	PCA	0.971	0.288	0.432	0.562	-9.56
	kPCA	0.870	0.444	0.547	0.651	-4.76
Нормальная смесь	PCA	0.583	0.571	0.547	0.567	0.14
	kPCA	0.651	0.573	0.575	0.539	-0.34
К-средних	PCA	0.611	0.769	0.653	0.585	0.41
	kPCA	0.732	0.513	0.571	0.580	-1.72
Байесовская нормальная смесь	PCA	0.613	0.810	0.671	0.564	0.52
	kPCA	0.728	0.513	0.572	0.595	-1.66

4:

3

Модель		Метрика качества				
		%	'		#>8B4	" 8BB >D=C
DBSCAN	PCA	0.401	0.930	0.539	0.352	3.74
	kPCA	0.460	0.834	0.551	0.352	2.68
Среднего сдвига	PCA	0.947	0.289	0.434	0.652	-13.00
	kPCA	0.821	0.472	0.530	0.727	-7.47
Нормальная смесь	PCA	0.434	0.607	0.476	0.535	2.74
	kPCA	0.492	0.616	0.495	0.536	1.68
К-средних	PCA	0.479	0.872	0.595	0.711	2.79
	kPCA	0.648	0.488	0.524	0.533	-1.37
Байесовская нормальная смесь	PCA	0.434	0.936	0.581	0.660	3.32
	kPCA	0.655	0.491	0.526	0.659	-1.53

5:

4

3.2

Наиболее точным оказывается метод среднего сдвига с линейным снижением размерности, при этом он является худшим по метрике полноты, которую удастся улучшить при использовании нелинейного снижения размерности.

Лучшим по *BCubed-полноте* оказывается DBSCAN.

По основной метрике *BCubed-F-мера* лучшим алгоритмом кластеризации является байесовская нормальная смесь.

Лучшим из рассмотренных способов снижения размерности почти всегда оказывается линейный метод главных компонент с отбором числа компонент по сохраненной дисперсии, но лучшим алгоритмом для построения шумового кластера оказывается метод среднего сдвига с нелинейным снижением размерности.

Стоит отметить, что более простой алгоритм К-средних показал наиболее близкое качество по метрикам кластеризации и выявления шума, как и байесовская нормальная смесь, при этом обучение модели К-средних проходит гораздо быстрее из-за меньшего количества параметров.

К применению в рамках задачи предлагается использование алгоритма кластеризации байесовской нормальной смесью.

4

Проведем дополнительный эксперимент, связанный с изучением вклада групп признаков в модель:

- признаки модели SBERT (*SBERT*);
- социально-демографические признаки (*соц.-дем.*);
- признаки, связанные с именованными сущностями и их тональностями (*NER*);

Так как при использовании метода главных компонент теряется интерпретируемость признаков, будем применять ее только для векторов SBERT, интерпретируемость которых и так отсутствует. В качестве алгоритмов кластеризации оставим только байесовскую нормальную смесь и К-средних как наилучшие по основной

метрике кластеризации, а в качестве метода снижения размерности оставим только линейный метод главных компонент с отбором числа компонент по сохраненной дисперсии. Способ подбора гиперпараметров и метрики качества будем использовать те же, рассматривать будем все темы совместно без срезов по числу мнений.

В качестве *базового* результата будем использовать кластеризацию обозначенным способом для всего векторного пространства, удаление соответствующих групп признаков будем обозначать *SBERT*, *соц.-дем.* и *NER*.

4.1

Модель		Метрика качества				
		%	'		#>8B4	" 8BB >D=C
Байесовская нормальная смесь	<i>базовый</i>	0.665	0.825	0.696	0.624	0.58
	<i>SBERT</i>	0.607 (-0.058)	0.819 (-0.006)	0.670 (-0.026)	0.600 (-0.025)	0.98 (+0.40)
	<i>соц.-дем.</i>	0.666 (+0.001)	0.823 (+0.004)	0.696 (+0.000)	0.625 (+0.001)	0.57 (-0.01)
	<i>NER</i>	0.661 (-0.004)	0.821 (-0.004)	0.691 (-0.005)	0.620 (-0.004)	0.54 (-0.04)
К-средних	<i>базовый</i>	0.685	0.764	0.680	0.640	0.24
	<i>SBERT</i>	0.642 (-0.043)	0.648 (-0.116)	0.609 (-0.071)	0.589 (-0.051)	0.13 (-0.11)
	<i>соц.-дем.</i>	0.677 (-0.008)	0.780 (+0.016)	0.686 (+0.006)	0.642 (+0.002)	0.44 (+0.20)
	<i>NER</i>	0.674 (-0.011)	0.771 (+0.007)	0.675 (-0.005)	0.640 (+0.000)	0.34 (+0.10)

6:

4.2

Признаки модели SBERT оказались наиболее важными для обоих алгоритмов кластеризации с точки зрения прироста по основной метрике кластеризации. Группа признаков, связанная с именованными сущностями и тональностями, также вносит положительный вклад в модель кластеризации. А вот группа социально-демографических признаков либо не влияет качество модели по *BCubed-F-мере*, как вышло для байесовской нормальной смеси, либо вовсе ухудшает модель, как вышло для алгоритма К-средних.

Заметим также, что кластеризация байесовской нормальной смесью в прошлом эксперименте лучше на 0.002 по *BCubed-F-мере*, чем кластеризация байесовской нормальной смесью в базовом варианте – лучше снижать размерность над всем признаковым пространством, а не только над пространством модели SBERT.

5

В рамках выпускной квалификационной работы удалось:

1. Реализовать алгоритм кластеризации для имеющихся признаков, позволяющий достигать качества **0.698** по метрике *BCubed-F-мера*;
2. Сравнить модели на способность выявлять шум в темах по предложенной метрике;
3. Предложить методику оценивания моделей кластеризации тем для поляризованных выборок;

- [1] E. Amigó и др. “A comparison of extrinsic clustering evaluation metrics based on formal constraints.” B: (2009).
- [2] D. Arthur и S. Vassilvitskii. “k-means++: The Advantages of Careful Seeding”. B: (2007).
- [3] D. Blei и M Jordan. “Variational inference for Dirichlet process mixtures”. B: (2006).
- [4] R. Cohen и D. Ruths. “Classifying Political Orientation on Twitter: It’s Not Easy!” B: (2013).
- [5] D. Comaniciu и P. Meer. “Mean Shift: A robust approach toward feature space analysis”. B: (2002).
- [6] Kareem Darwish и др. “News Consumption in Time of Conflict: 2021 Palestinian-Israel War as an Example”. B: (2021).
- [7] A.P. Dempster, N.M. Laird и D.B. Rubin. “Maximum Likelihood from Incomplete Data via the EM Algorithm”. B: (1977).
- [8] Subhabrata Dutta и др. “Semi-supervised Stance Detection of Tweets Via Distant Network Supervision”. B: (2022).
- [9] M. Ester и др. “A Density-Based Algorithm for Discovering Clusters in Large Spatial Databases with Noise”. B: (1996).
- [10] Tiziano Fagni и Stefano Cresci. “Fine-Grained Prediction of Political Leaning on Social Media with Unsupervised Deep Learning”. B: (2022).
- [11] T. Fawcett. “An Introduction to ROC Analysis”. B: (2006).
- [12] C.J. Filmore. *Some problems for case grammar*. 1971.
- [13] Margherita Gambini и др. “Tweets2Stance: Users stance detection exploiting Zero-Shot Learning Algorithms on Tweets”. B: (2022).
- [14] S. Geman и D. Geman. “Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Image”. B: (1984).

- [15] M. Korobov. “Morphological Analyzer and Generator for Russian and Ukrainian Languages”. English. В: *Analysis of Images, Social Networks and Texts*. Под ред. М. Khachay и др. Т. 542. Communications in Computer and Information Science. Springer International Publishing, 2015, с. 320—332. ISBN: 978-3-319-26122-5. DOI: [10.1007/978-3-319-26123-2_31](https://doi.org/10.1007/978-3-319-26123-2_31). URL: http://dx.doi.org/10.1007/978-3-319-26123-2_31.
- [16] Jinning Li и др. “Unsupervised Belief Representation Learning in Polarized Networks with Information-Theoretic Variational Graph Auto-Encoders”. В: (2022).
- [17] Jun. S. Liu. “The Collapsed Gibbs Sampler in Bayesian Computations with Applications to a Gene Regulation Problem”. В: (1994).
- [18] Jun Lu. “A survey on Bayesian inference for Gaussian mixture model”. В: (2021).
- [19] K. Pearson. “On Lines and Planes of Closest Fit to Systems of Points in Space”. В: (1901).
- [20] F. Pedregosa и др. “Scikit-learn: Machine Learning in Python”. В: *Journal of Machine Learning Research* 12 (2011), с. 2825—2830.
- [21] Peter J. Rousseeuw. “Silhouettes: a Graphical Aid to the Interpretation and Validation of Cluster Analysis”. В: (1987).
- [22] Lloyd S. “Least square quantization in PCM’s.” В: (1957).
- [23] SberDevices. “BERT large model multitask (cased) for Sentence Embeddings in Russian language”. В: (2021). URL: https://huggingface.co/sberbank-ai/sbert_large_mt_nlu_ru.
- [24] B. Schölkopf, A. Smola и K. Müller. “Nonlinear Component Analysis as a Kernel Eigenvalue Problem”. В: (1998).
- [25] К.В. Воронцов. *Вероятностное тематическое моделирование: теория, модели, алгоритмы и проект VigARTM*. 2021.
- [26] Д. Фельдман и К.В. Воронцов. “Комбинирование фактов, семантических ролей и тональных слов в генеративной модели для поиска мнений”. В: (2020).

Выявление поляризации текстов новостей внутри тем.

На странице задания представлена группа текстов, которые принадлежат в основном какой-то одной теме. Ваша задача: прочитать все тексты, выделить мнения-полюсы, которые по Вашему мнению присутствуют в данной группе новостей, и определить, к какому из выделенных вами мнений относится каждый текст. Под полюсом тут подразумевается попытка авторов текста отразить некоторую смещенную точку зрения на определенное событие, например:

1.
 - Полюс 1: «Мнение журналистов ВВС»
 - «Журналист ВВС заявил, что британский эсминец намеренно нарушил границы РФ в Черном море»;
 - «Корреспондент ВВС раскрыл правду об инциденте с «Дефендером» в Черном море»;
 - Полюс 2: «Взгляд Российских властей»
 - «Россия заявила о рисках проведения США и их союзниками учений в Черном море»;
 - «Россия призвала США отказаться от военных маневров в Черном море»;
 - «Сенатор от Крыма оправдала действия российских бомбардировщиков в отношении британского эсминца»;
2.
 - Полюс 1: «Произвол/негатив»
 - «В вытрезвитель теперь смогут забрать даже из дома, к тому же у клиента будет произведен досмотр вещей»;
 - Полюс 2: «Нейтральное/констатация факта»;
 - «МВД РФ согласовало правила помещения россиян в медицинские вытрезвители»;

