

А.Ю. Хоменко, Т.В. Романова

АЛГОРИТМ АВТОМАТИЧЕСКОЙ АТРИБУЦИИ ПИСЬМЕННОГО ТЕКСТА В ЛИНГВОКРИМИНАЛИСТИКЕ

В рамках исследования была сделана попытка поиска алгоритма, с одной стороны, основанного на интеграции интерпретационного подхода к анализу языкового знака и методов математической статистики и математического моделирования, а с другой - подпадающего автоматизации.

Целью работы явилось определение, могут ли выбранные методы математической статистики и стилистического анализа успешно применяться в судебно-авторствовании (определить автора текста чаще всего требуется в спорах, связанных с нарушением авторских и смежных прав (ст. 146 УК РФ) можно ли на их основе создать универсальную безошибочную методику атрибуции текста *любого объема* и можно ли автоматизировать процесс атрибуции текста на современном этапе развития российской судебно-экспертной практики).

Соответственно, ставилась задача - разработать алгоритм, включающий как методы интерпретационного анализа, так и методы математической статистики и теории вероятности и позволяющий как можно более точно в автоматическом режиме определять автора письменного текста.

Материалы и методика исследования. Материалом для исследования послужили художественные тексты заведомо известных авторов, поскольку целью работы является определение того, *развешивать ли методика стилистического анализа для текстов различных стилей и объемов*.

1) Пестовая выборка (ТВ) - выборка, на основе которой строилась исходная модель. ТВ представляла собой тексты С.Д. Довлатова, размещенные в Национальном корпусе русского языка (URL: <http://www.nkorpus.ru/>), за исключением текста «Наши» (1983 г.) (ЭТ). Этот текст рассматривался как экспериментальный, то есть текст, у которого якобы не определен автор. Таким образом, объем ТВ - 330769 слов; 2) ЭТ (экспериментальный текст) - текст С. Довлатова «Наши» (1983 г.), текст, автор которого якобы неизвестен. Объем 21230 слов.

В основу методики анализа было положено исследование Е.С. Родионовой¹, где описываются методы стилистического анализа, квантитативная лингвистика и теории распознавания образов. Методика Е.С. Родионовой была совмещена с методикой анализа языковой личности по Ю.Н. Караулову², мето-

дикой квантитативного анализа незамысловатых и стилистически немаркированных лексем А.Н. Баранова и некоторыми постулатами теории вероятности.

В конечном итоге был разработан алгоритм со следующими этапами:

1. Построение атрибуционных типов об авторстве спорных текстов (текст текстов, автора якобы неизвестны);

2. Анализ языковой личности (ЯЛ). Отбор параметров для будущей математической модели. Результатом стали 35 выявленных характеристик ЯЛ С.Д. Довлатова на трех уровнях языковой личности: 1) *вербально-семантический уровень*: местоимения «я», «мы», «ты», «они». Эти характеристики выделены, исходя из наличия соотношения в прозе Довлатова проблемы субъекта повествования и субъекта действия, то есть автора и героя. Использование сочинительных союзов «и», «а», «но» в начале предложения; 2) *лексико-синтаксический уровень*: «лексемы, маркирующие отношение к действительности». Эти лексемы отражают в том числе аксиологические оценки; б) «*объемный*», «*белый*», «*светлый*», «*русский*», «*русский*», «*русский*», «*молочный*», «*начальный*» - лексемы, имплицитно маркирующие отношение к действительности и создающие её образ; в) «*город*», «*человек*», «*детство*», «*родина*» - лексемы, вербализующие значимые для Довлатова концептуальные сущности. Образ города присутствует во многих произведениях Довлатова, являясь символом определённого типа сознания. Чемодан является «эксплицитно модальности допущения, неуверенности. Ну, тусня, давно, но-жасвай, так, бы, видно; б) эксплицитно модальности удивления: *неужели, разве, ах, ах*; в) эксплицитно модальности ограничения: *только, лишь, почти; з) эксплицитно модальности возражения: *вот-таки**.

3. Квантитативные и стилистические преобразования данных, полученных в результате того, сколько раз параметр реализуется в каждой выборке; 2) Определены среднестатистические частоты по формуле (1):

$$\bar{x} = \frac{1}{n} \sum_{i=1}^k x_i = \frac{1}{n} (x_1 + \dots + x_k)$$
 формула 1, - где x_i - i -й элемент выборки, n - объем выборки; 3) Определение отклонения выборочных частот от среднестатистической частоты по формуле (2):

$$s = \sqrt{\frac{n-1}{n} \sum_{i=1}^k x_i^2 - \bar{x}^2} \quad \text{формула 2, - где } \sigma^2 \text{ — дисперсия; } x_i \text{ — } i\text{-й элемент выборки; } n \text{ — объем выборки; } \bar{x} \text{ — среднее арифметическое выборки (среднестатистическая частота); 4) Поиск вероятной ошибки в определении средней частоты по формуле (3) (для } \alpha = 0,2 \text{ и вероятности } 0,8 \text{ при } (n - 1) \text{ степеней свободы (} 35 - 1 = 34 \text{)): } t = 1,3070$$

$$L = \frac{k}{\sqrt{n}}$$
 формула 3, - где t - табличный коэффициент (t-критерий Стьюдента); σ - среднестатистическое отклонение; k - объем выборки. Для ТВ ошибка составляет

¹ Баранова А.Н. Введение в прикладную лингвистику: Учебное пособие, 2001. URL: <http://ebooks.vsu.ru/docs/books/index-418241.html>

¹ Родионова Е.С. Лингвистические методы атрибуции и алгоритм интерпретации произведений проблемы «Молшер - Корсилье». Автореферат диссертации на соискание степени кандидата филологических наук, 2008. URL: <http://eprints.uva.nl/archive/00000553/00000553.pdf>

² Караулов Ю.Н. Русский язык и языковая личность. М.: Наука, 1987

ставляет 0,002272751. Для ЭТ - 0,008969957. Ошибка для ТВ и ЭТ незначительна, но тем не менее, она учитывается во матрицах данных; 5) Определение релевантных параметров для конечных моделей. Определяются по критерию Стюарта (4). Уровень значимости $\alpha = 0,2$. Критическое значение - в таблице перемены уровня степеней свободы (количества параметров - 1) и вероятия стп 0,8.

$$t = \frac{\bar{x}_1 - \bar{x}_2 - \mu_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \text{ формула 4. - где } \bar{x}_1, \bar{x}_2 - \text{ средние арифметические; } s_1^2, s_2^2 - \text{ стандартные отклонение; } n_1, n_2 - \text{ объемы выборок.}$$

По результатам исследования выделены следующие релевантные для моделей параметры: модель ТВ и ЭТ: *эрусья, сочинительный союз "но" в начале предложения, муть, эрусьяный разе, ах, бейлы, нурусья*. Релевантными для построения моделей в настоящей работе считаются параметры, числовые показатели которых наиболее близки к табличному значению (критерия (1, 3070)). Важно, что ни одно значение параметра не превысило значение t-статистики. Для настоящего исследования это означает, что полученные результаты будут иметь точность менее изначально заявленной, то есть менее 80%.

IV. Переход от реальных объектов к их математическим моделям, то есть описание выделенных в ходе предыдущего анализа параметров с помощью условной сигнатуры. Формирование матриц данных; введение двух моделей: модели текстов-образов, описывающей некоторое количество языковой единицы заведомо известного автора, и модели спурного текста, описывающей некоторые закономерности языковой личности автора. Сравнения моделей осуществляется коэффициентом корреляции между однородными параметрами (5). Этот коэффициент показывает, насколько близки две модели. Чем ближе значение этого коэффициента к 1, тем более близки модели качествами отношения.

$$r_{xy} = \frac{\sum_{i=1}^n x_i y_i}{\sqrt{\sum_{i=1}^n x_i^2 \sum_{i=1}^n y_i^2}} \text{ формула 5. - где}$$

r_{xy} - коэффициент корреляции между значениями матриц ТВ и ЭТ равно 0,783448911306154.

Как видим, коэффициент корреляции достаточно чуть ниже 80%, тем не менее, даже уровень приближенности в 78%, на наш взгляд, можно считать значительным.

Таким образом, подтвердилось следующая апробированная гипотеза: **Н₀** - автор ТВ и ЭТ - одно лицо, то есть автор **ТВ и ЭТ - С.Д. Довгань** (по закону граничности: если автор **ТВ - С.Д. Довгань**, а автор **ЭТ и ТВ** - одно лицо, то автор **ЭТ - тоже С.Д. Довгань**).

Дальнейшие разработки в выбранной сфере позволили выработать рекомендации для улучшения работы алгоритма.

Так, для использования методики в условиях реальной действительности и для текстов большого объема можно дать следующие рекомендации:

- ▶ число параметров для идентификации автора по письменному речевому прогнозноному должно быть не менее 45 – 50 единиц;
- ▶ параметры должны представлять собой обширные синтаксические и рфонологические классы, например, большинство матричных респондентов субъективной модальности (автоные слова, модальные частицы, междометия, конструкторы с инвариантным представлением);
- ▶ отбор параметров должен происходить на основе более глубокого анализа языковой личности автора текста-образца, причем в большем объеме можно на мотивационном уровне (возможно, также на вербально-матричном). Для применения указанной методики без внесения в саму методику изменений можно дать следующие рекомендации:
- ▶ объемы текста-образца и спурного текста должны быть близки и не должны превышать 10 страниц печатного текста каждый;
- ▶ тексты должны быть схожи с точки зрения функциональных стилей;
- ▶ методу можно дополнить вычислением из двух текстов (эптонного текста, то есть сравнительного образца, и спурного текста), так называемых разноминимных текстов.

Е.Ю. Салаева, В.М. Бухаров

ПРОСОДИКА СПОНТАННОЙ РЕЧИ: МЕТОДЫ ФОРМИРОВАНИЯ ЭКСПЕРИМЕНТАЛЬНОГО КОРПУСА И ЕГО АНАЛИЗА

Настоящий этап развития лингвистики характеризуется повышенным вниманием к звуковой стороне речи с функциональной точки зрения. Просодика, обеспечивая интонационное выразительность, играет важную роль для осуществления коммуникативной функции, так как именно с помощью просодических средств мы передаем то, что не можем передать на сегментном уровне. Благодаря исследованиям в данной области стало возможным создание систем просодического распознавания голоса (по характеристикам темпа, ритма, мелодического оформления речи). Такие системы широко внедряются в нашу жизнь. Например, для оценки качества работы операторов контакт-центров, совершенствующие подобная систем требуют дальнейших исследований в области просодики. А благодаря развитию речевой лингвистики, начиная с 60х годов прошлого века, в нашем распоряжении появляются речевые корпуса как база для лингвистических исследований. Речевой корпус – это «структурированное множество речевых фрагментов, которое обеспечено программными средствами доступа к отдельным элементам корпуса». Для одних целей создается достаточным использованием уже существующих корпусов, когда для других задач необходимо создание отдельного экспериментального корпуса.

Эти и другие разработки в области речевого телекода и мультимедийной фонетики принадлежат исследователям компании ООО «Лингвистика Технологии». Подготовлены специалистами ГИТ Липовича также в сотрудничестве с Сергеем РЯД, АЛЕКСАНДРОМ ТЕЛЕЦКОМ, АЛЕКСАНДРОМ МАСИ и др. ООО «Лингвистика Технологии» работает на основе технологии интеллектуальной «Специализированный корпус СЛД ОР» <http://www.speechcorp.ru> (лицензия 04.02.2014)