

УДК 81'33

DOI: 10.18384/2310-712X-2015-4-71-77

**Литвинова Т.А., Литвинова О.А.***Воронежский государственный педагогический университет***ИССЛЕДОВАНИЕ ЛИНГВИСТИЧЕСКИХ ХАРАКТЕРИСТИК ТЕКСТОВ,  
СОДЕРЖАЩИХ НАМЕРЕННО ИСКАЖЁННУЮ ИНФОРМАЦИЮ,  
С ПОМОЩЬЮ ПРОГРАММЫ *LINGUISTIC INQUIRY AND WORD COUNT***

*Аннотация.* В статье представлены результаты исследования, направленного на выявление статистически значимых различий между «правдивыми» текстами и текстами, содержащими намеренно искажённую информацию. Исследование проводилось на материале специально созданного корпуса текстов – «Банка лжи», содержащего один «правдивый» и один «ложный» текст на одну и ту же тему от каждого автора, а также метаданные в виде информации об авторах (пол, возраст, данные психологического тестирования и т.д.). Все тексты были обработаны при помощи программы Linguistic Inquiry and Word Count. Статистический анализ позволил выявить значимые различия между «правдивыми» и «ложными» текстами по ряду лингвистических характеристик.

*Ключевые слова:* ложь, лингвистика лжи, корпус текстов, компьютерная лингвистика, ложный текст, языковые маркеры лжи.

**T. Litvinova, O. Litvinova***Voronezh State Pedagogical University***A STUDY OF LINGUISTIC FEATURES OF DECEPTIVE TEXTS WITH  
THE USE OF THE PROGRAM *LINGUISTIC INQUIRY AND WORD COUNT***

*Abstract.* The results of the study dealing with the identification of statistically significant differences between «truthful» texts and those containing intentionally deceptive information are presented. The study was performed using a specially designed text corpus «Deception Bank» with a «truthful» and a «deceptive» text by the same author on the same topic as well as the metadata including the information about the author (gender, age, psychological testing results, etc.). All the texts were processed using Linguistic Inquiry Word Count. The statistical analysis revealed that a wide range of linguistic characteristics differed considerably in «truthful» and «deceptive» texts.

*Key words:* deception, text-based deception detection, text corpus, computer linguistics, deceptive text, linguistic markers of deception.

Несмотря на то, что проблема выявления в речевой продукции лжи, под

© Литвинова Т.А., Литвинова О.А., 2015.

Исследование выполнено при поддержке гранта РГНФ 15-34-01221 «Детекция лжи в письменном тексте: корпусное исследование».

которой мы понимаем «намеренно созданный продукт мыслительной деятельности человека, искажённо (полностью или частично) отражающий действительность» [8, с. 346], имеет не только теоретическую, но и очевидную практи-

Таблица 1

**Коэффициенты вариации параметров (в процентах),  
рассчитанные на основе текстов «Банка лжи»**

Параметр	% слов длинной более 6 букв	% строєвых слов	% местоиме- ний	% личных местоиме- ний	% мест. «я» и косв. форм	% мест. «мы» и косв. форм	% мест. «она, он» и косв. форм	% мест. «они» и косв. форм
Коэф. вариации	8	7	17	17	25	41	42	40
Параметр	% наречий	% предлогов	% союзов	% отрицаний	% числительных	% слов группы «Эмоции»	% слов группы «Положительные эмоции»	% слов группы «Негативные эмоции»
Коэф. вариации	36	10	25	47	59	23	26	55
Параметр	% слов группы «Мыслительные процессы»	% слов группы «Интуиция»	% слов группы «Причинно- следственные отношения»	% слов группы «Несоответствие»	% слов группы «Попытка»	% слов группы «Уверенность»	% слов группы «Помеха»	% слов группы «Объединение»
Коэф. вариации	13	34	40	47	47	37	54	27
Параметр	% слов группы «Исключение»	% слов группы «Восприятие»	% слов группы «Зрительное восприятие»	% слов группы «Слуховое восприятие»	% слов группы «Ощущение»	% слов группы «Пространство»	% слов группы «Время»	% знаков препинания
Коэф. вариации	38	37	45	47	46	14	18	8

ческую значимость [1; 7], она относится к числу малоизученных в российской лингвистике. Учёными признаётся тот факт, что на различных языковых уровнях возможно «обнаружение типологических свойств лжи», а сама ложь, в ко-

торой «существует особая грамматика, особое использование лексики, особые правила словоупотребления и синтаксиса», определяется как «особый код в тексте, или как “язык в языке”» [9, с. 35]. Однако системного описания типоло-

Таблица 2

**Усреднённые значения параметров (в процентах от общего числа слов)  
для «правдивых» (П) и «ложных» (Л) текстов**

	% слов длинной более 6 букв	% строковых слов	% местоиме- ний	% личных местоиме- ний	% мест. «я» и косв. форм	% мест. «мы» и косв. форм	% мест. «она, он» и косв. форм	% мест. «они» и косв. форм
П	30,06	40,50	13,11	10,32	5,61	1,83	2,42	1,99
Л	29,66	40,88	13,43	10,95	5,93	1,92	2,54	1,83
	% наречий	% предлогов	% союзов	% отрицаний	% числительных	% слов группы «Эмоции»	% слов группы «Положитель- ные эмоции»	% слов группы «Негативные эмоции»
П	2,27	13,33	4,43	2,07	0,52	6,78	5,08	1,22
Л	2,17	13,24	4,47	2,12	0,59	6,82	5,28	1,11
	% слов группы «Мыслительные процессы»	% слов группы «Интуиция»	% слов группы «Причины- следственные отношения»	% слов группы «Несоответствие»	% слов группы «Попытка»	% слов группы «Уверенность»	% слов группы «Помеха»	% слов группы «Объединение»
П	13,77	2,64	1,67	0,84	1,51	1,84	0,49	5,08
Л	14,04	2,88	1,76	0,77	1,55	1,80	0,46	5,25
	% слов группы «Исключение»	% слов группы «Восприятие»	% слов группы «Зрительное восприятие»	% слов группы «Слуховое восприятие»	% слов группы «Ощущение»	% слов группы «Пространство»	% слов группы «Время»	% знаков препинания
П	1,49	2,61	0,80	0,91	0,70	11,08	9,66	21,84
Л	1,50	2,29	0,72	1,00	0,52	11,19	9,22	20,34

гических свойств письменных текстов на русском языке, содержащих ложную информацию, до настоящего времени выполнено не было.

В зарубежной науке наблюдается повышенный интерес к проблеме выявления лжи в речевой продукции, и в письменном тексте в частности. Такие исследования имеют в основном междисциплинарный характер. Лингвисты, психологи и специалисты по математическому моделированию объединяют свои усилия, чтобы выявить статистически значимые различия между прав-

дивыми текстами и текстами, содержащими ложь (в основном исследования проводятся на материале английского языка), а затем на основе найденных различий построить математические модели, позволяющие с определённой вероятностью установить наличие или отсутствие в тексте намеренно искажённой информации [5; 6].

**Целью** нашей работы является поиск регулярных статистически значимых различий между правдивыми и «ложными» текстами на русском языке. Насколько нам известно, данная

Таблица 3

**Изменение значений параметров в «ложных» текстах  
относительно «правдивых» (в процентах)**

	% слов длинной более 6 букв	% строєвых слов	% местоиме- ний	% личных местоиме- ний	% мест. «Я» и косв. форм	% мест. «мы» и косв. форм	% мест. «она, он» и косв. форм	% мест. «они» и косв. форм
Изменение параметра в «ложном» тексте относительно «правдивого», %	-1,34	0,95	2,43	6,06	5,74	4,76	5,22	-7,63
	% наречий	% предлогов	% союзов	% отрицаний	% числительных	% слов группы «Эмоции»	% слов группы «Положительные эмоции»	% слов группы «Негативные эмоции»
	-4,54	-0,62	0,80	2,45	13,16	0,62	3,94	-8,41
	% слов группы «Мыслительные процессы»	% слов группы «Интуиция»	% слов группы «Причинно- следственные отношения»	% слов группы «Несоответствие»	% слов группы «Попытка»	% слов группы «Уверенность»	% слов группы «Помеха»	% слов группы «Объединение»
	1,97	9,07	4,94	-7,96	2,50	-1,57	-6,57	3,45
	% слов группы «Исключение»	% слов группы «Восприятие»	% слов группы «Зрительное восприятие»	% слов группы «Слуховое восприятие»	% слов группы «Ощущение»	% слов группы «Пространство»	% слов группы «Время»	% знаков препинания
	0,99	-11,95	-10,21	10,45	-25,47	1,05	-4,65	-6,86

задача на материале русского языка решается впервые.

Для решения поставленной задачи необходим прежде всего исследовательский корпус, содержащий как «правдивые», так и «ложные» тексты. Составление такого рода корпусов – самостоятельная научная проблема [11; 12]. До настоящего времени такой корпус текстов на русском языке отсутствовал, а работы, посвящённые изучению типологических свойств лжи в речи, проводились на материале художественных текстов [5; 6]. В связи с этим в 2014 г. нами был начат сбор «Банка лжи» в составе корпуса текстов *Personality* [4]. В настоящее время

«Банк лжи» содержит 228 текстов от 114 респондентов, в дальнейшем планируется расширение корпуса. Помимо правдивого и «ложного» текста от каждого автора на одну и ту же тему («Как я провёл вчерашний день»), подкорпус, как и все тексты корпуса *Personality*, содержит метаданные в виде информации об их авторах – пол, возраст, данные психологического тестирования. Таким образом, аннотированный «Банк лжи» позволит в дальнейшем произвести учёт индивидуальных особенностей автора при продуцировании им «ложного» текста. На наш взгляд, без учёта этих данных невозможно построение объективной

методики выявления в тексте намеренно искажённой информации, хотя проблема учёта индивидуальных различий при анализе речевой продукции на предмет наличия в ней ложной информации только начинает привлекать внимание учёных [6].

Очевидно, что ни одна языковая категория не может служить маркером лжи сама по себе; важны лишь частотные характеристики некоторых параметров текста в их совокупности, составляющие своеобразный «лингвистический профиль лжи». Для анализа текста по целому ряду лингвистических параметров целесообразно привлечение программных средств. В нашем исследовании для обработки текстов мы использовали программу *Linguistic Inquiry and Word Count*, которая позволяет получать числовые значения параметров текста (процент слов определённых грамматических, лексико-семантических категорий и т.д. от общего числа слов в тексте) [13] и активно используется в зарубежных исследованиях проблемы взаимосвязи языка и личности, а также для выявления различий между «правдивыми» и «ложными» текстами [10]. Анализ был проведён в соответствии с 32 параметрами, по которым предположительно могут различаться «правдивые» и «ложные» тексты.

Задача выявления маркеров лжи в письменном тексте решалась нами с использованием методов математической статистики в три этапа<sup>1</sup>. На первом этапе мы определяли коэффициент вариации каждого параметра

в «правдивых» и «ложных» текстах, чтобы понять, значения каких лингвистических характеристик остаются стабильными в текстах одного автора, а каких – варьируются. Для этого мы рассчитали отклонение значений каждого параметра текста от его средней величины для конкретного автора. Далее мы усреднили отклонения каждого параметра по текстам всех авторов и определили коэффициент вариации параметров, что позволило нам оценить меру разбросанности значений параметров текста и понять, насколько она велика относительно их среднего значения. Расчёты выполнялись при помощи статистического пакета *SPSS*.

Статистический анализ (табл. 1) показал, что рассчитанный коэффициент вариации для выбранных параметров текста лежит в достаточно широких пределах. Принимая во внимание тот факт, что значение коэффициента вариации менее 33 % свидетельствует об однородности совокупности данных; параметры с таким значением коэффициента вариации мы считаем устойчивыми в текстах одного автора.

Для того чтобы понять, как по абсолютной величине изменяются параметры текста, мы рассчитали усреднённые значения каждого параметра (табл. 2). В таблице 3 приведено относительное изменение каждого параметра в «ложных» текстах относительно «правдивых».

На основе анализа полученных данных мы предлагаем ввести весьма удачный, на наш взгляд, маркер, который может быть использован для сравнения «правдивых» и «ложных» текстов, – отношение процентного содержания наречий в тексте к процентному содержанию личных место-

<sup>1</sup> Авторы выражают глубокую благодарность доктору физико-математических наук П.В. Середину за помощь в статистической обработке материала.

имений. Выбор в пользу этого маркера был сделан нами по следующим причинам. Во-первых, средние значения параметров «процентное содержание наречий в тексте» и «процентное содержание личных местоимений в тексте» выше единицы (см. табл. 2), т.е. оба эти параметра являются **частотными**. Во-вторых, как видно из таблицы 3, изменения параметров в случае «ложного» текста происходят во взаимно **противоположных** направлениях, т.е. в среднем в «ложных» текстах по сравнению с «правдивыми» текстами одного и того же автора наречий становится меньше, а личных местоимений – больше. В-третьих, исходя из анализа стабильности / вариативности параметров текста (см. табл. 1), коэффициент вариации для наречий в два раза выше, чем для личных местоимений, и превышает величину 33 %, что свидетельствует о **нестабильности** этого параметра, которая может объясняться, в числе прочего, наличием / отсутствием в тексте намеренно искажённой информации.

Выбор в пользу **соотношения** двух указанных параметров сделан ещё и потому, что появляется возможность сравнивать тексты разного объёма и исключить влияние личностных особенностей автора на значения параметров текста. Ранее нами было установлено, что частотность употребления в тексте слов тех или иных частей речи, а именно местоимений и служебных слов, коррелирует с некоторыми параметрами личности [2; 3]. Соотношения параметров, каждый из которых зависит и от характеристик личности, исключают такое влияние.

Исходя из полученных нами данных, для текстов «Банка лжи» мы рас-

считали усреднённое значение выбранного нами параметра, а далее – ошибку в его определении. Мы получили, что для «правдивых» текстов значение выбранного параметра должно быть более **0,21±0,02**, а для «ложных» – менее **0,19±0,02**.

Для проверки эффективности предложенного нами параметра – маркера лжи – мы использовали два подкорпуса, составленных на основе корпуса текстов *Personality*. Первый подкорпус, который мы использовали в качестве тестового корпуса «правдивых» текстов (312 текстов), представлен текстами – описаниями картин. Второй использованный нами тестовый корпус был составлен из текстов, авторы которых пытались убедить воображаемых работодателей выбрать их для занятия вакантной должности, намеренно искажая информацию о себе (корпус «ложных» текстов, 64 текста). Для всех текстов были рассчитаны значения выбранного нами маркера лжи. Точность определения «правдивых» текстов (для первого подкорпуса) составила ~71 %, «ложных» (для второго подкорпуса) текстов – ~72 %.

Таким образом, нами с использованием программы *Linguistic Inquiry and Word Count* на материале специально подготовленных корпусов текстов было проведено исследование, доказавшее существование статистически значимых различий между «правдивыми» и «ложными» текстами одного автора. Предложенный нами параметр позволяет определить факт наличия / отсутствия в тексте намеренно искажённой информации с вероятностью ~71–72 %, причём этот результат был получен на разных корпусах текстов, что доказывает правильность выбран-

ного нами подхода. В дальнейшем необходима проверка эффективности предложенного параметра на новом корпусном материале. Исследование также показало, что поиск эффективных маркеров лжи должен проводиться с обязательным учётом уровня стабильности / вариативности параметра в текстах одного автора.

#### ЛИТЕРАТУРА:

1. Леонтьев А.А. Прикладная психолингвистика речевого общения и массовой коммуникации. М.: Смысл, 2008. 272 с.
2. Литвинова Т.А. Формально-грамматические корреляты личностных особенностей автора письменного текста // Филологические науки. Вопросы теории и практики. 2013. № 12 (30). Ч. 1. С. 132–135.
3. Литвинова Т.А. Языковые корреляты личностных особенностей автора письменного текста: алгоритм исследования // В мире научных открытий. Серия: Проблемы науки и образования. 2012. № 9.3(33). С. 236–255.
4. Литвинова Т.А. и др. Корпусные исследования письменной речи в решении задач судебного автороведения / Литвинова Т.А., Диброва Е.В., Литвинова О.А., Рыжкова Е.С. // Филологические науки. Вопросы теории и практики. 2015. № 8. С. 56–69.
5. Литвинова Т.А., Литвинова О.А. Исследование текста на предмет наличия в нём намеренно искажённой информации: проблемы и перспективы // Известия ВГПУ. Серии: «Педагогические науки»; «Гуманитарные науки». 2015. № 2 (267). С. 189–192.
6. Литвинова Т.А., Середин П.В. Поиск признаков лжи в письменном тексте: современные методы и подходы // В мире науки и искусства: вопросы филологии, искусствоведения и культурологии: материалы XXIV Международной заочной научно-практической конференции [10 июня 2013 года] / СибАК. Новосибирск: СибАК, 2013. С. 126–133.
7. Орлова Н.В. Проблема правды / лжи в судебной лингвистической экспертизе. [Электронный ресурс]. URL: [http://siberia-expert.com/publ/satti/stati/problema\\_pravdy\\_lzhi\\_v\\_sudebnoj\\_lingvisticheskoj\\_ekspertize\\_n\\_v\\_orlova/4-1-0-204](http://siberia-expert.com/publ/satti/stati/problema_pravdy_lzhi_v_sudebnoj_lingvisticheskoj_ekspertize_n_v_orlova/4-1-0-204) (дата обращения: 20.05.2015).
8. Потапова Р.К., Потапов В.В. Язык. Речь. Личность. М.: Языки славянской культуры, 2006. 496 с.
9. Степанов Ю.С. Альтернативный мир, дискурс, факт и принцип причинности // Язык и наука XX века: сб. ст. М.: РГГУ, 1995. С. 35–73.
10. Newman M. et al. Lying words: Predicting deception from linguistic style / Newman M., Pennebaker J., Berry D., Richards J. // Personality and Social Psychology Bulletin. 2003. № 29. P. 665–675.
11. Rubin V., Conroy N. J. The art of creating an informative data collection for automated deception detection: A corpus of truths and lies // Proceedings of the American Society for Information Science and Technology. 2012. Vol. 49. N 1. DOI: 10.1002/meet.14504901045 [Электронный ресурс]. URL: [http://www.researchgate.net/publication/259540155\\_The\\_art\\_of\\_creating\\_an\\_informative\\_data\\_collection\\_for\\_automated\\_deception\\_detection\\_A\\_corpus\\_of\\_truths\\_and\\_lies](http://www.researchgate.net/publication/259540155_The_art_of_creating_an_informative_data_collection_for_automated_deception_detection_A_corpus_of_truths_and_lies) (дата обращения: 20 мая 2015).
12. Salvetti F. Detecting Deception in Text: A Corpus-Driven Approach: Ph.D. dissertation. University of Colorado at Boulder, 2012. 205 p.
13. The Development and Psychometric Properties of LIWC2007 / J.W. Pennebaker [et al.] [Электронный ресурс]. URL: [http://homepage.psy.utexas.edu/homepage/faculty/pennebaker/reprints/liwc2007\\_languagemanual.pdf](http://homepage.psy.utexas.edu/homepage/faculty/pennebaker/reprints/liwc2007_languagemanual.pdf) (дата обращения: 19.05.2015).