

УДК 81'322::811.222.8::519.25
ББК 81.1::Ш152.131.2-923::В172

МУАЙЯНКУНИИ МУАЛЛИФИ АСАРҲОИ ТАЪРИХИЮ-СИЁСӢ БО ВОСИТАИ БИГРАММАҲОИ РАМЗӢ

Қосимов А.А. – номзади илмҳои техникаӣ, омӯзгори калон, кафедраи барномарезӣ ва низомҳои иттилоотӣ, Донишқадаи политехникаи Донишгоҳи техникаи Тоҷикистон ба номи академик М.С. Осимӣ,
abdunabi.kbtut@mail.ru

Қаюмов А.А. – магистрант соли дуюми ихтисоси информатика, Академияи миллии илмҳои Тоҷикистон Маркази илмии Хуҷанд,
zohirjon96@gmail.com

Чакида. Муқаррар карда шуд, ки басомади воҳурии биграммаҳои забони тоҷикӣ дар асарҳои таърихӣ - сиёсӣ ин муайянкунандаи муаллиф аст. Барои муайянкунии муаллифи матн бо басомади воҳурии ҳарфҳои биграмма ва имконияти истифодабарии таснифгари Усмонов З.Ҷ. таҳқиқот гузаронида шуд. Роҳҳо, мақсад ва имкониятҳои таснифгари Усмонов З.Ҷ. оварда мешаванд. Симои рақамӣ ва фазои ченаки асарҳои таърихӣ-сиёсӣ сохта шудаанд. Бо мақсади ягонагии эҷодиёти муаллиф, қимати ченаки муқоисакунанда муқаррар карда мешаванд, ки дар асоси он ягонагии синфҳои асарҳо муайян мешаванд. ӯ-таснифгари бузургии тасодуфии фосиладор, ки самаранокии баландро дар муайян кардани муаллифи асарҳои шоирону нависандагони форсу тоҷик нишон дода буданд, барои муайян кардани муаллифи асарҳои таърихӣ-сиёсӣ тафтиш карда мешаванд.

Калидвожаҳо: забони тоҷикӣ, асарҳои таърихӣ-сиёсӣ, биграмма, таснифгар, басомади воҳурий.

Натиҷаҳои, ки олимони дар ин ҷо ба даст оварданд, фақат барои асарҳои бадеӣ ва иқтисодӣ буданд. Аз ҳамин нуқтаи назар дар ин мақола маълумот оид ба таҳқиқоти муайян кардани муаллифи асарҳои таърихӣ-сиёсӣ бо воситаи басомади воҳурии биграммаҳои рамзӣ оварда мешавад. Масъалаи шинохти матн, дар асл, аз он рӯзе, ки хат пайдо шуд, ба миён омадааст. Дар муддати тӯлони, он танҳо бо яке аз “рӯя”-ҳои худ, ки дар он зарурати муайян кардани иҷрокунандаи асари навишташуда ифода меёфт, мавриди таҳлил ва баррасӣ қарор додбиграма мешуд. Минбаъд, пас аз лаҳзаи ихтироъ шудани “китобчопкунӣ”, дар ин

масъала, рӯяи нави муҳими дигари он – зарурати муайян кардани муаллифи маводи чопӣ ба миён омад. Дар марҳилаи кунунӣ маҳз ҳамин “рӯя”-и дуюм, мазмуни мундариҷаи асосии ин масъаларо ташкил медиҳад.

Дар кори зерин ба сифати инструменти таҳқиқотшаванда таснифгари Усмонов З.Ҷ. санҷида мешавад, [3-5].

Маълумот оид ба коллексияи матнҳо. Асарҳои муаллифҳои зерин гирифта шуданд ва дар дохили қафс шакли кӯтоҳ карда шудаи онҳо, ки барои ҷойгиркунии ҷадвалҳои 1 ва 2-и поёни лозим буд, оварда шудаанд: А. Дюма (АД) “Граф монте кристо 1” (ГМК1) ва “Граф монте кристо 2” (ГМК2) [6, 7], А. Хуш (АХ) “Шӯриши Усмон 1” (ШУ1) ва “Шӯриши Усмон 2” (ШУ2) [8-10], М. Шакурӣ (МШ) “Садри Бухоро” (СБ) ва “Хуросон аст инчо” (Х) [11, 12], С. Айни (СА) “Ёддоштҳо 1” (Ё1) ва “Ёддоштҳо 2” (Ё2) [13, 14], Ҷ. Икромӣ (ҶИ) “Дувоздаҳ дарвозаи Бухоро” (ДБ) ва “Он чи аз сар гузашт” (О) [15, 16]. Ба омӯзиши ин масъала ҳамагӣ 5-муаллиф ва 10-асар гирифта шуд.

Таснифгари матнҳо. Ба сифати тавсири рақамии асарҳои таърихӣ-сиёсӣ басомади вохӯрии биграммаҳои рамзӣ дида баромада мешавад. Барои муайянкунии муаллиф як метод – таснифгари матнии Усмонов З.Ҷ. истифода бурда шуд. Моҳияти тавсифи ин метод дар татбиқ ба масъалаҳои илми забоншиносӣ дода мешавад, [3].

Бигзор T_1 ва T_2 – ду матне бошанд, ки қонуни тақсимооти биграммаҳои рамзии онҳо ба намуди ҷадвал дода шуда бошад

$$\begin{aligned} T_i &: 1 \dots k \dots m \\ P^{(i)} &: p_1^{(i)} \dots p_k^{(i)} \dots p_m^{(i)}, \end{aligned} \quad (1)$$

ки дар ин ҷо

$$\sum_{k=1}^m p_k^{(i)} = 1 \quad \text{аст.}$$

Дар ин ифодаҳо k ($k = \overline{1, m}$) - рақами тартибии биграммаи k -юм дар алифбои биграмма, $p_k^{(i)}$ - басомади нисбии вохӯрии биграммаи k -юм дар матни T_i , $i = 1, 2$ мебошад. Он гоҳ масофаи байни T_1 ва T_2 бо формулаи зерин муайян карда мешавад

$$\rho(T_1, T_2) = \sqrt{\frac{m}{2}} \max_s \left| \sum_{k=1}^s (p_k^{(1)} - p_k^{(2)}) \right|, \quad (2)$$

дар ин ҷо $s = \overline{1, m}$.

Бигзор γ - ягон адади мусбат бошад, матнҳои T_1 ва T_2 γ -якҷинса номида мешаванд, агар

$$\rho(T_1, T_2) \leq \gamma. \quad (3)$$

ва γ -ғайриякҷинса номида мешаванд, агар

$$\rho(T_1, T_2) > \gamma \text{ бошад.} \quad (4)$$

Фарз мекунем, ки коллексияи матнҳо T ба зермаҷмӯъҳои $T^{(j)}$, $j = \overline{1, n}$ тақсим шудааст. Барои қиммати қайдшудаи γ адади \aleph^0 -суммаи ҷуфтҳои якҷинсаи матн, ки ба зермаҷмӯъҳои $T^{(j)}$, $j = \overline{1, n}$, тааллуқ доранд ва адади \aleph^H -суммаи γ -ҷуфтҳои ғайриякҷинса, ки ба зермаҷмӯъҳои гуногун тааллуқ дорад, ҳисоб карда мешавад. Нисбати

$$\eta = \frac{\aleph^0 + \aleph^H}{N}, \quad (5)$$

ки дар ин ҷо N -шумораи умумии ҷуфти матнҳо дар коллексияи T аст, барои қиммати дода шудаи γ самаранокии татбиқи модели математикии (1) – (4) ба таври автоматӣ тақсимкунии коллексияи T ба зерқисми $T^{(j)}$ -ро тавсир мекунад. Дар мақолаҳои [4, 5], барои ҳисоб кардани қимати оптималии γ^{opt} , ки барои он самаранокии максималии η барои коллексияи $T = \{T^{(j)}\}$ дастрас мегардад, пешниҳод гардид.

Натиҷаҳо. Алгоритми дар боло зикр шударо истифода бурда, комплекси барномаҳо тартиб дода шуданд ва дар аввал басомади вохӯрии биграммаҳои рамзӣ бе ва бо ҳисобгирии фосола дар алоҳидагӣ ҳисоб карда, баъдан масофаи байни асарҳо бо формулаи (2) ҳисоб карда шуданд, натиҷаҳо дар ҷадвалҳои 1 ва 2 оварда шудаанд. Аз натиҷаҳои ба даст омада чунин қонуниятро бояд ҷудо кард, ки ду асари як муаллиф якҷинсаанд ва ду асари муаллифашон гуногун ғайриякҷинсаанд.

Чадвали 1

Қиммати γ^{omm} барои биграмма бе ҳисобгирии фосила

АД		АХ		МШ		СА		ҚИ	
ГМК1	ГМК2	ШУ1	ШУ2	СБ	Х	Ё1	Ё2	ДБ	О
0,2009									
0,6662	0,5189								
0,5392	0,4037	0,3346							
1,0775	0,9086	0,5243	0,5713						
1,2595	1,0893	0,6072	0,7271	0,3460					
0,4132	0,4545	0,7709	0,7056	1,1627	1,3366				
0,3545	0,3225	0,6334	0,5647	1,0520	1,1492	0,2366			
0,4808	0,4015	0,3344	0,4129	0,7288	0,9068	0,5041	0,3901		
0,4242	0,3144	0,3357	0,3016	0,7034	0,8951	0,5116	0,3791	0,3400	

Бояд қайд кард, ки дар ҳар ду ин чадвал дар диагонали асосӣ маълумот оид ба муносибати байни асарҳои як муаллиф, аммо дар дигар ячейкаҳо маълумот оид ба муносибати байни асарҳои муаллифашон гуногун оварда шудаанд.

Барои муайянкунии муаллиф диапазони қиммати мувофиқи γ бо воситаи биграмма бе ҳисобгирии фосила баробари $\gamma = [0,0301; 0,0346)$ шуд. Дар ин ҳолат бо нобаробарии зерин

$$\rho(T_1, T_2) \leq [0,0301; 0,0346) \quad (6)$$

- якҷинсагии ҷуфти асарҳо, аммо бо муқобили нобаробарии,

$$\rho(T_1, T_2) > [0,0301; 0,0346), \quad (7)$$

- бо ғайриҷинсагии асарҳо мувофиқат мекунад. Ин қоидаро ба қатори ададҳои чадвали 1 татбиқ намоем, нобаробарии (6) дар 16 ячейка ғайриҷинсагии асарҳо намешавад, аммо нобаробарии (7) фақат дар 2 ячейка ғайриҷинсагии асарҳо намешавад. Ба ҳолати зерин, таъсири метод бо формулаи (5) ҳисоб карда шуд, ки баробари $\eta = 91\%$ аст.

Ҷадвали 2

Қиммати γ^{opt} барои биграмма бо ҳисобгирии фосила

АД		АХ		МШ		СА		ЧИ	
ГМК1	ГМК2	ШУ1	ШУ2	СБ	Х	Ё1	Ё2	ДБ	О
0,0503									
0,1434	0,1610								
0,2027	0,1741	0,0597							
0,3242	0,3665	0,2263	0,2190						
0,3718	0,4148	0,2934	0,2688	0,0843					
0,2587	0,2374	0,3413	0,3499	0,5257	0,5935				
0,2278	0,1931	0,3142	0,3208	0,4961	0,5679	0,0927			
0,1265	0,1055	0,2164	0,2305	0,4346	0,4771	0,1648	0,1281		
0,1236	0,1153	0,1975	0,2041	0,3717	0,4439	0,2037	0,1751	0,1245	

Таҳлили ҷадвали 2 нишон медиҳад, ки бо воситаи биграмма бо ҳисобгирии фосила, таснифгар ҳиссиёти баландтарро доро аст. Барои муайянкунии муаллиф диапазони қиммати мувофиқи γ бо воситаи биграмма бо ҳисобгирии фосила баробари $\gamma = [0,0927; 0,1245)$ шуд. Дар ин ҳолат бо нобаробарии зерин

$$\rho(T_1, T_2) \leq [0,0927; 0,1245) \quad (8)$$

- якҷинсагии чуфти асарҳо, аммо бо муқобили нобаробари,

$$\rho(T_1, T_2) > [0,0927; 0,1245), \quad (9)$$

- бо ғайриякҷинсагии асарҳо мувофиқат мекунад. Ин қоидаро ба қатори ададҳои ҷадвали 2 татбиқ намоем, нобаробарии (8) дар 10 ячейка рияо намешавад, аммо нобаробарии (9) бошад, фақат дар 2 ячейка рияо намешавад. Ба ҳолати зерин, таъсирнокии метод бо формулаи (5) ҳисоб

карда шуд, ки баробари $\eta = 93\%$ аст.

Хулоса. Аз маълумотҳои ҳангоми таҳқиқот ба даст омада, ба чунин хулосаҳо омадан мумкин аст, ки

- биграммаи рамзӣ дар масъалаи муайянкунии муаллифи матни таърихӣ-сиёсӣ ба сифати тавсифҳои миқдорӣ комилан қобили қабул мебошанд;

- ба ҳисобгирии фосила дар биграммаҳо саҳеҳии таснифотро баланд мебардорад;

- таснифгари Усмонов З.Ҷ. (1) – (5) дараҷаи кифоя калони муайянкунии муаллифи матнҳои таърихӣ-сиёсиро нишон медиҳад.

ЛИТЕРАТУРА

1. Усмонов З.Д., Косимов А.А. Частотность биграмм в таджикской литературе.– Доклады Академии наук Республики Таджикистан, 2016, т.59, № 1-2, с. 28-32.

2. Романов А.С., Шелупанов А.А., Мещеряков Р.В. Разработка и исследование математических моделей, методик и программных средств информационных процессов при идентификации автора текста.– Томск: -В-Спектр, 2011, 188 с.

3. Усмонов З.Д. N-граммы в распознавании однородных текстов.– Материалы 20 научно-практического семинара "Новые информационные технологии в автоматизированных системах".– Москва, 2017, с. 52-54.

4. Усмонов З.Д. Классификатор дискретных случайных величин.– Доклады Академии наук Республики Таджикистан. 2017, т.60, № 7-8, с. 291-300.

5. Усмонов З.Д. Алгоритм настройки кластеризатора дискретных случайных величин.– Доклады Академии наук Республики Таджикистан, 2017, т.60, № 9, с. 392-397.

6. Александр Дюма Граф монте кристо. 1844, 520 с.

7. Алӣ Хуш Шӯриши Усмон . Душанбе, 2008, 128 с.

8. Муҳаммадҷон Шакурӣ Садри Бухоро. Душанбе, 2005, 330 с.

9. Муҳаммадҷон Шакурӣ Хуросон аст инчо.

10. Садриддин. Айнӣ Ёддоштҳо. 1990, 352 с.

11. Ҷалол. Икромӣ Дувоздаҳ дарвозаи Бухоро. 1969, 474 с.

12. Ҷалол. Икромӣ Он чи аз сар гузашт. Душанбе, 2009, 360 с.

ИДЕНТИФИКАЦИЯ АВТОРОВ ИСТОРИЧЕСКОЙ ПОЛИТИЧЕСКОЙ ПРОИЗВЕДЕНИЙ С ПОМОЩЬЮ СИМВОЛНЫХ БИГРАММ

Косимов А.А. – кандидат технических наук, старший преподаватель, кафедра программирования и информационных технологий, Политехнический институт Таджикского технического университета имени академика М.С. Осими, abdunabi_kbtut@mail.ru

Каюмов А.А. – магистрант, Академия наук Республики Таджикистан, Худжандский научный центр, zohirjon96@gmail.com

Аннотация. Устанавливается, что распределение частотности биграмм в исторических-политических произведениях таджикского языка является идентификатором авторства. Исследованы возможности классификатора З.Д.Усманова распознавать автора текста по частотности буквенных биграмм. Сконструированы цифровой портрет и метрическое пространство произведений. В предположении уникальности авторского творчества устанавливаются пороговые значения метрики, на основе которых определяются классы “однородных” произведений. -классификатор дискретных случайных величин, подтвердивший высокую эффективность при идентификации авторства текстовых фрагментов в произведениях классической и современной поэзии, а также в современной прозе таджикского языка, тестируется на предмет приспособляемости к распознаванию авторства в исторической -политической произведениях.

Ключевые слова: таджикский язык, исторический-политический произведения, биграмма, классификатор, частотность.

IDENTIFICATION OF AUTHORS OF A ECONOMICS-POLITICAL TEXT BY MEANS OF A SYMBOL BIGRAMS

Kosimov A.A. – Candidate of technical Sciences, Department of Programming and Information Technologies, Polytechnic Institute of Tajik Technical University, abdunabi_kbtut@mail.ru

Kayumov A.A. – master student, National Academy of sciences of Tajikistan, Khujand scientific center, zohirjon96@gmail.com

Annotation. *It is established that the distribution of the frequency of bigrams in the historical and political works of the tajik language is the identifier of authorship. The possibilities of the classifier Z.D.Usmanov to recognize the author of the text by the frequency of letter bigrams are investigated. A digital portrait and a metric space of works are constructed. Assuming the uniqueness of the author's creativity, the threshold values of the metrics are established, on the basis of which the classes of "homogeneous" works are determined. The γ -classifier of discrete random variables, which confirmed high efficiency in identifying the authorship of text fragments in works of classical and modern poetry, as well as in modern prose of the Tajik language, is tested for adaptability to the recognition of authorship in historical-political text.*

Key words: *tajik language, historical-political text, bigram, classifier, frequenc.*