

ئىنتوناتسىيە قاتلاملىرىنى ئاپتوماتىك ئايرىشقا قارشا ئېلىپ بېرىلغان ئۇيغۇر تىلى سۆز تۈركۈملىرىنى ئاپتوماتىك بەلگىلەش تېخنىكىسى تەتقىقاتى

نۇربىيە تاھىر¹، دىلمۇرات تۇرسۇن^{1*}، ئەسقەر روزى²

1(شىنجاڭ ئۇنىۋېرسىتېتى ئۇچۇر پەن-تېخنىكىسى ۋە ئىنژىنىرلىقى ئىنستىتۇتى، ئۈرۈمچى
(830046

2(شىنجاڭ ئۇنىۋېرسىتېتى ماتېماتىكا ۋە سېستىما پەن تېخنىكىسى
ئىنستىتۇتى، ئۈرۈمچى (830046

ئۈزۈندە: بۇ ماقالە تىل بېرىكتۈرۈش سېستېمىسىنىڭ تىكىست ئانالىز قىلىش مودىلىدىكى ئىنتوناتسىيە چىگرىسىنى ئاپتوماتىك ئايرىش تېخنىكا بۆلىكى ئارقا كۆرۈش قىلىنىپ، ئۇيغۇر تىلىدىكى سۆز تۈركۈملىرىنى ئاپتوماتىك بەلگىلەش تېخنىكىسى نوقتىلىق تەتقىق قىلىندى. يەنى، ئالدى بىلەن قوللىنىش ساھەسىدىكى ئالاھىدىلىكىگە ئاساسەن سۆز تۈركۈملىرىنىڭ سانى ۋە بىكىتىش پىرىنسىپى بەلگىلىنىپ جۈملىلەر تاللاپ چىقىلدى ھەمدە ئۇنىڭغا قارشا ئادەم كۈچى ئارقىلىق ھەر بىر سۆزنىڭ جۈملە مۇھىتىدىكى سۆز تۈركۈمى بەلگىلەندى؛ ئاندىن سىتاتىستىكا ئۇسۇلى ئارقىلىق سۆز تۈركۈملىرىنىڭ نىسپى جەدۋىلى ۋە سۆز تۈركۈملىرىنىڭ ماسلىشىش جەدۋىلى ھاسىل قىلىندى؛ ئاخىردا HMM مودىلىنىڭ ئىككى ئىلمىنىڭ ئۇسۇلى ئارقىلىق ئۇيغۇر تىلىنىڭ سۆز تۈركۈمىنى ئاپتوماتىك بەلگىلەش تېخنىكىسى ئەمەلگە ئاشۇرۇلدى. تەجرىبە جەريانىدا قوللىنىلغان ئالگورىتمىنىڭ ئۈنۈمى تەكشۈرۈش ئۈچۈن 10000 جۈملە مەشقلەندۈرۈلۈپ، توپلام ئىچى ۋە سىرتىدىن ئايرىم-ئايرىم ھالدا 500 دانەدىن جۈملە تاللاپ ئېلىندى ھەمدە ئۇلارنىڭ سۆز تۈركۈملىرى ئاپتوماتىك بەلگىلىنىپ ئۇنىڭ توغرىلۇق دەرىجىسى تەكشۈرۈلدى. تەجرىبە نەتىجىسى ئانالىز ئۇسۇلىنىڭ قوللىنىش قىممىتى بارلىقىنى دەلىللەپ بەردى.

ھالقىلىق سۆزلەر: ئۇيغۇر تىلى؛ ئىنتوناتسىيە قاتلىمى؛ سۆز تۈركۈمىنى بەلگىلەش؛ تىل بېرىكتۈرۈش؛ تىكىست ئانالىزى

Research on a POS tagging method for detecting the boundary of prosody in Uyghur Sentences

Nurbiye Tahir¹, Dilmurat Tursun^{1*}, A skar Rozi²

¹ (College of Information Science and Engineering, Xinjiang University, Urumqi 830046, China)

² (College of Mathematics and System Science, Xinjiang University, Urumqi 830046, China)

[Abstract]: As a critical issue in the text analysis part of Uyghur synthesis system, an automatic POS (part of speech) tagging method for estimation the boundary of different prosody levels in text sentences is mainly discussed in this paper. At first, according to the characteristics of specific application field, the category of parts used in POS tagging and the rules are confirmed, candidate sentences are selected and are tagged by manually, then POS probability list tables and relative occurring position information tables of POS are achieved by statistics, and at last, the automatic POS tagging method is implemented through adopting bigram model on the basis of HMM model. In order to approve the validity of the method presented in this paper, the large scale text corpus is used in which 10000 sentences are selected for training and 500 sentences are used for testing. Test results show that the method adopted in this paper for Uyghur POS tagging is feasible, appropriate and valid.

[Keywords]: Uyghur; prosodic layer; POS tagging; speech synthesis; text analysis

1. كىرىش سۆز

2. سۆز تۈركۈمى نىسپەت جەدۋىلى

يېقىنقى يىللاردىن بۇيان، ئاز سانلىق مىللەت تىل-يېزىق ئۇچۇر بىر تەرەپ قىلىش ساھەسىدىكى يادرولقۇ تېخنىكىلارنىڭ بىرى بولغان-ئادەم ۋە ماشىنا ئالاقىسى تېخنىكىسى، ئۇيغۇر تىل بېرىكتۈرۈش تەتقىقاتىدا ئاساس بولۇش نەزىرىيە تەتقىقاتى ۋە تېخنىكا قوللۇنۇلۇش قاتارلىق تەرەپلەردە كۆرۈنەرلىك ئىلگىرىلەشلەرگە ئېرىشتى. ئەمما بېرىكتۈرۈلگەندىن كىيىنكى ئاۋازنىڭ تەبئىيلىكى يەنىلا كۆزلىگەن نىشاندىن خېلىلا يىراق بولماقتا. بېرىكتۈرۈلگەن ئاۋازنىڭ تەبئىيلىكىنى تېخىمۇ يۇقىرى كۆتۈرۈش ئۈچۈن ئاۋاز ئامبىرىنى قايتىدىن قۇرۇش، مەسىللەر توپلىمىنى ئەلالاشتۇرۇش، تىل بېرىكتۈرۈش جەريانىدىكى ئالدىنقى بۆلەك مودىل قۇرلىشىنى ئىنچىكە لايىھىلەش، ئىنتوناتسىيەلىك بۆلەكلەر چىگرىسىنى ئالدىن پەرەز قىلىپ بەلگىلەش، ئۇرغۇنى ئالدىن پەرەز قىلىش، سۆز

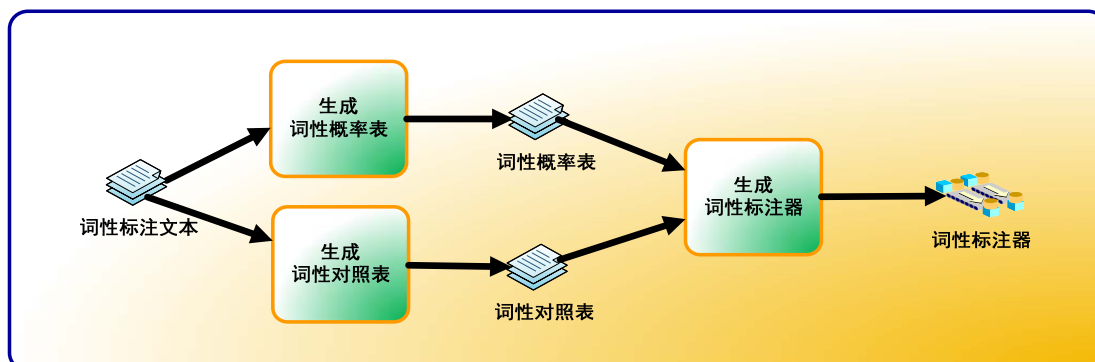
تۈركۈمىنى ئالدىن بەلگىلەش ئۇسۇلىنى ئەلاشتۇرۇش قاتارلىق تەرەپلەردە يەنىمۇ بىر قەدەم ئىلگىرلىگەن ئاساستا تەتقىق قىلىشقا ۋە ئىزدىنىشكە توغرا كېلىدۇ. شۇڭلاشقا ئېنتوناتسىيە بۆلەكلىرىنىڭ چىگرىسىنى ئالدىن پەرەز قىلىپ بەلگىلەش بۇ ساھەدىكى قىززىق نوقتىغا، شۇنداقلا بېرىكتۈرۈلگەن ئاۋازنىڭ تەبئىيلىك دەرىجىسىنى يۇقىرى كۆتۈرۈشتىكى مۇھىم ھالقىلىق مەسىلىگە ئايلاندى.

ئادەتتە تىل بېرىكتۈرۈش سېستىمىسى تۆۋەندىكىدەك تۆت قىسىمنى ئۆز ئىچىگە ئالىدۇ: تېكىست ئانالىز قىلىش مودىلى، ئېنتوناتسىيە ھاسىل قىلىش مودىلى، ئاۋاز ئىلمى مودىلى، ئاۋاز ئاساسى بىرلىكى ۋە ئاۋاز پارامىتىرى ئامبىرى مودىلىدىن ئىبارەت.^[1] ئېنتوناتسىيە قەۋەتلىرىنى ئايرىش جەريانىدا، كىرگۈزۈلگەن تېكىستنىڭ خاس ئۇچۇرلىرى بولۇش بىلەن بىرگە تېكىستتىكى ھەر قايسى سۆزلەرنىڭ سۆز تۈركۈمى "ئاپتوماتىك سۆز تۈركۈمى بەلگىلىگۈچ" ئارقىلىق ئېرىشىلگەن بولۇشى كېرەك.^[2] ئۇيغۇرچە تىل بېرىكتۈرۈش سېستىمىسىدىكى ئېنتوناتسىيەلىك تەركىبلەرنىڭ چىگرىسىنى ئاپتوماتىك ئايرىش بولسا تېكىست ئانالىز قىلىش مودىلىدىكى مۇھىم بىر ھالقا بولۇپ ھېساپلىنىدۇ. ئېنتوناتسىيەلىك تەركىبلەرنىڭ چىگرىسىنى ئايرىشنىڭ توغۇرلۇق دەرىجىسىنى ئاشۇرۇش ئۈچۈن ئەڭ ئۈنۈملۈك بولغان سۆز تۈركۈمىنى ئاپتوماتىك بەلگىلەش ئۇسۇلىنى تاللاپ قوللىنىش كېرەك. بۇ خىزمەتنىڭ ۋەزىپىسى بولسا، كونكرىت تىل مۇھىتىدە ھەر بىر سۆزگە ئېنىق بىر سۆز تۈركۈمىنى بەلگىلەپ چىقىشتىن ئىبارەت. بۇ جەزىم دەرىخى مودىلى بىلەن ماسلىشىپ ئېنتوناتسىيەلىك تەركىبلەرنىڭ چىگرىسىنى ئاپتوماتىك ئايرىش جەريانىدا مۇھىم رول ئوينايدۇ. ئۇيغۇرتىلىدىكى ئېنتوناتسىيەلىك تەركىبلەرنىڭ چىگرىسىنى ئاپتوماتىك ھالدا توغرا ئايرىپ چىقىش بولسا، ئەڭ ئاخىردا تىل بېرىكتۈرۈش سېستىمىسىدا بېرىكتۈرۈلگەن ئاۋازنىڭ تەبئىيلىكىنى يۇقىرى كۆتۈرۈشكە خىزمەت قىلىدۇ. شۇڭلاشقا ئېرىشمەكچى بولغان نەتىجە ئۇيغۇرچە تىل بېرىكتۈرۈش تەتقىقاتىدا ناھايىتى مۇھىم ئورۇندا تۇرىدۇ.

بۇ خىزمەت جەريانىدا بەلگىلەنگەن سۆز تۈركۈمىنىڭ توغرىلىقىغا كاپالەتلىك قىلىش ئۈچۈن ئالدى بىلەن 10610 جۈملەگە قاراش خاس مۇھىتتىكى سۆز تۈركۈمىنى بەلگىلەش خىزمىتى ئېلىپ بېرىلدى. بۇ سانلىق مەلۇماتلار ئۇيغۇرتىلى ئېنتوناتسىيە تەتقىقات خىزمىتىنىڭ چۇڭقۇرلاپ قانات يايدۇرۇلۇشى ئۈچۈن ناھايىتى ياخشى ئاساس سېلىپ بەردى، بۇ يەنە ئۆز نۆۋىتىدە بېرىكتۈرۈلگەندىن كېيىنكى ئاۋازنىڭ تەبئىيلىك دەرىجىسىنى ئاشۇرۇش، پۈتكۈل سېستىمىنىڭ تېكىست ئانالىز قىلىش مودىلىدىكى تېخنىكىلىق خىزمەتلەرنىڭ چۇڭقۇرلاپ قانات يايدۇرۇلۇشى ئۈچۈنمۇ ئاساس بولغۇسى! بۇ خىزمەت جەريانىدا 10610 جۈملىدىن 10000 جۈملە تاللىۋېلىنىپ مەشقلەندۈرۈلدى، توپلام ئىچى ۋە سىرتىدىن ئايرىم-ئايرىم ھالدا 500 جۈملە تاللىۋېلىنىپ سىناق قىلىندى. بۇ تۈردىكى تەتقىقاتلار ئاۋاز بىر تەرەپ قىلىش تەتقىقاتىنىڭ قوللىنىش ساھەلىرىدە يەنى: تىل بېرىكتۈرۈش، تىلنى تونۇش، ئاۋاز تەتقىقاتى، تىل تەتقىقاتى قاتارلىق ساھەلەردە ئىلمى ۋە ئەمەلىي قىممىتى بار. بۇنىڭ يەنە ئۇيغۇرتىلى ھەتتا پۈتكۈل ئالتاي تىللىرى سېستىمىسىدىكى تىللارنىڭ ئاۋاز تەتقىقاتى قوللىنىش ۋە ئېچىش جەريانىدىمۇ خېلى يۇقىرى پايدىلىنىش قىممىتى بار.

2. سۆز تۈركۈمىنى ئاپتوماتىك بەلگىلەش سېستىمىسىنىڭ مودىللاشقان لايىھەسى

بۇ خىزمەتتىكى سۆز تۈركۈمىنى ئاپتوماتىك بەلگىلەش سېستىمىسى: سۆز تۈركۈمىنى بەلگىلەش تېكىست مودىلى، سۆز تۈركۈمى سىتاسىتىكا ئۇچۇرلىرىغا ئېرىشىش مودىلى، سۆز تۈركۈمىنى ئاپتوماتىك بەلگىلىگۈچ مودىلىدىن تەركىپ تاپقان. 1-رەسىمدە كۆرسىتىلگەندەك:



1-رەسىم سۆز تۈركۈمىنى ئاپتوماتىك بەلگىلەش سېستىمىسىنىڭ قۇرۇلمىسى

2.1 سۆز تۈركۈمىنى بەلگىلەش تېكىست مودىلى

2.1.1 سۆز تۈركۈمى ۋە ئۇنىڭ سانى

سۆز تۈركۈمىنى بەلگىلەش ئېنگىلىز تىلىدا (Part-of-Speech tagging OR POS tagging) دېيىلىدۇ. بۇ ئامباردىكى ھەر بىر سۆزگە ماس بولغان تىپ يەنى سۆزلەرنىڭ تەۋەلىك گۇرۇپپىسىنى بەلگىلەپ چىقىش جەريانىدىن ئىبارەت. گەرچە تەبئىي تىل-ئالاقە جەريانىدىكى بەلگىلەر تېخىمۇ كۆپ مەنەلەرگە ئىگە بولۇشتەك ئالاھىدىلىككە ئىگە بولسىمۇ، ئەمما تەبئىي تىلنىڭ ئالاقە جەريانىدىكى ئالاھىدىلىكلىرىنى بەلگىلەش جەريانى ماشىنا تىلىدىكى سۆزلەرنىڭ ئەسلىگە قايتۇرۇلۇش (tokenization) جەريانى بىلەن ئوخشىشىپ كېتىدۇ^[2].

سۆز تۈركۈمىنىڭ تىل ئۇچۇرلىرىنى بىر تەرەپ قىلىشتىكى ئەھمىيىتى-سۆز ۋە ئۇنىڭ ئەتراپىدىكى بۆلەكلەرگە مۇناسىۋەتلىك كۆپلىگەن ئۇچۇرلارنى تەمىنلەپ بېرىلگەنلىكىدە. شۇڭلاشقا سۆز تۈركۈمىنى ئاپتوماتىك بەلگىلەش ئاۋاز پەرىقلەندۈرۈش، تەبئىي تىلنى تەھلىل قىلىش ۋە ئۇچۇر ئىزدەش قاتارلىق ساھەلەردە بارغانسېرى كۆپ قوللىنىلماقتا. سۆز تۈركۈملىرى بەلگىلىنىپ بولغان ماتېرىيال ئامبىرى تىل تەتقىقاتىدا ناھايىتى مۇھىم ئەھمىيەتكە ئىگە. تىل بېرىكتۈرۈش سېستېمىسىدا ئەگەر بىر سۆزنىڭ سۆز تۈركۈمىنى بىلىسەك تېخىمۇ تەبئىي بولغان ئاۋازغا ئېرىشەلەيمىز. ئۇيغۇر تىلىدىكى سۆز تۈركۈملىرىنىڭ تۈرلىرى بولسا سۆز يىلتىزلىرىنىڭ گىرامماتىكىلىق ئالاھىدىلىكلەر بويىچە بۆلۈنىشىدىن ئىبارەت. گىرامماتىكىلىق ئالاھىدىلىكلەر بولسا: شەكىل ئۆزگۈرۈش، سۆز بىلەن سۆزنىڭ بېرىكىش ئىقتىدارى، سۆزنىڭ جۈملىدىكى ئىقتىدارى ۋە سۆزنىڭ قۇرۇلۇش ئالاھىدىلىكى قاتارلىقلاردىن ئىبارەت. ئوخشاش بىر تۈركۈمگە تەۋە بولغان سۆزلەر ئوخشاش گىرامماتىكىلىق ئالاھىدىلىككە ئىگە بولىدۇ.

ئىلگىرىكى ئۇيغۇرچە تىل بېرىكتۈرۈش سېستېمىسىنىڭ سۆز تۈركۈمىنى بەلگىلەش جەريانىدا 16 تۈرلۈك سۆز تىپى قوللىنىلغان^[3]. بۇ خىزمەت جەريانىدا بېيجىڭ ئۇنۋېرسىتېتىنىڭ سۆز تۈركۈمىنى بەلگىلەش قائىدىسى ۋە ئىلگىرى قوللانغان 16 خىل سۆز تۈركۈمى ئاساسىدا، ئېنتوناسىيەگە كۆرسىتىدىغان تەسىرى بىر قەدەر چوڭ دەپ قارالغان 41 خىل سۆز تۈركۈمى دەسلەپكى قەدەمدە بېكىتىلدى. سىتاتىستىكا خىزمىتى ئاخىرلاشقاندا سۆز تۈركۈمىنىڭ سانىنى مۇقۇملاشتۇرۇشقا قاراپ يەنە بىر قېتىم جەزىملەشتۈرۈش خىزمىتى ئىشلىنىپ ئەڭ ئاخىرىدا 37 خىل سۆز تۈركۈمى مۇقۇملاشتۇرۇلدى. بۇ سۆز تۈركۈملىرىنىڭ ئىپادىلىنىشى تۆۋەندىكى 1-جەدۋەلدە قىسمى كۆرسىتىلگەندەك:

1-جەدۋەل سۆز تۈركۈملىرى جەدۋىلى

文标识	英文标识	维文标识	序号	中文标识	英语标识	维语标识	序号
不定代词	rb	بەلگىسىز ئالماش	20	名词	n	ئىسىم	1
反身代词	ro	ئۆزۈك ئالماش	21	人名	nr	كىشى ئىسىملىرى	2
副词	d	رەۋىشى	22	地名	ns	يەر-جاي نامى	3
情态副词	dh	ھالەت رەۋىشى	23	位置名词	nt	ئورگان تەشكىلات نامى	4
耐间副词	dw	ۋاقىت رەۋىشى	24	行业术语	nz	كەسپى ئاتالغۇلار	5
位置副词	do	ئورۇن رەۋىشى	25	形容词	a	سۈپەت	6

2.1.2 تىكىست لايىھەسى

تىكىست لايىھەسىنىڭ مەقسىتى: ئامالنىڭ بارىچە سان جەھەتتىن ئاز بولغان تىكىست ئارقىلىق تېخىمۇ كۆپ سۆز تۈركۈمى ھادىسىنى ئۆز ئىچىگە ئېلىش. بۇنداق بولغاندا ھەر خىل جۈملە شەكىللىرى ۋە ھەر خىل سۆز تۈركۈملىرىنىڭ كۆرىلىش نىسپىتىگە كاپالەتلىك قىلغىلى بولىدۇ.

جۈملە تاللىغاندا جۈملە ئۇزۇنلىقى نەزەرگە ئېلىنىپ، جۈملە شەكلى بويىچە تۈرگە ئايرىپ تاللاندى. ئالدى بىلەن ئەسلىدىكى تېكىستنى بوغۇم، سۆز، سۆز بېرىكىمىسى ۋە جۈملىلەرگە پارچىلاپ ئاندىن Greedy ئالگورىتمى ئارقىلىق زور كۈلەمدىكى ماتېرىيال ئامبىرىدىن ۋەكىللىك خاراكتېرىگە ئىگە بولغان جۈملىلەر توپلىمى ھاسىل قىلىندى^[3]. تىكىستنى تاللاپ بولغاندىن كىيىن ئۇنۇڭغا قاراپ ئادەم كۈچى ئارقىلىق ئۆزگەرتىش، تاللاش ۋە شاللاش

خىزمىتى تەكرار ئېلىپ بېرىلىپ كەمتۈك جايلارنى تولۇقلاندى.

2.1.3 ئادەم كۈچى بىلەن سۆز تۈركۈمىنى بەلگىلەش

ئۇيغۇرتىلى يېپىشقاقلىققا ئىگە بىر تىل بولۇپ، شەكىل ئۆزگۈرۈشى ناھايىتى مۇرەككەپ. ئۇنىڭ ئۈستىگە بىر قانچە سۆز تۈركۈمىگە تەۋە بولىدىغان سۆزلەر بىر قەدەر كۆپ، قوللىنىلىش نىسپىتىمۇ بىر قەدەر يۇقىرى [4]. بۇ ماقالىدىكى سۆز تۈركۈمىنى بەلگىلەش ئاساسى-شۇ سۆزنى جۈملە مۇھىتىدە تەھلىل قىلىپ ئۇنىڭ سۆز تۈركۈمىنى بەلگىلەش. بۇ سەۋەپلىك ئادەم كۈچى ئارقىلىق بەلگىلىگەندە بەزىدە جۈملىدىكى شۇ سۆزنىڭ ئۆزىلا ئويلىشىشقا توغرا كەلسە بەزىدە سۆز تۈمۈرى، سۆز تۈرمۈرىنىڭ ئۆزگۈرۈشى، سۆز تۈمۈرى ۋە قوشۇمچىلارنىڭ ماسلىشىشى، ھەر خىل سۆز تۈرلىگۈچى قوشۇمچىلار ۋە سۆز ياسىغۇچى قوشۇمچىلارنىڭ ئۆز ئارا ماسلىشىشى قاتارلىق ئامىللارنى ئويلىشىشقا توغرا كېلىدۇ، شۇڭلاشقا بۇنىڭغا قارىتا بىرلىككە كەلگەن ئۆلچەم بېكىتىش قىيىن بولىدۇ. بەلگىلىگەنگەن سۆز تۈركۈمىنىڭ توغرىلىقىغا كاپالەتلىك قىلىش ئۈچۈن، بۇ گورۇپپىدىكى ئەزالار سۆزلەرنىڭ تۈركۈملەرگە بۆلىنىشى ۋە مۇناسىۋەتلىك بىلىملەرگە پىششىق بولغانلار ۋە ئۇزۇن يىللىق ئۇيغۇرتىلى گىرامماتىكا ئوقۇتۇش تەتقىقاتى بىلەن شۇغۇللانغان ئوقۇتقۇچىلاردىن تەشكىللەندى.

ئادەم كۈچى ئارقىلىق تەييارلانغان بۇ خىل ھۆججەتنىڭ مەزمۇنى تۆۋەندە كۆرسىتىلگەندەك:

1. ئۇيان do/بۇيان do/ئۆتۈشۈۋاتقان av/شۇنچە rk/كۆپ/ئادەملەر/ئارىسىدا do/شەيخ/نەنىلا/ئۇنىڭ ra/ئارقىسىدىن n/جىمغىنا/قاراپ dv/تۇراتتى v/.
2. قىزمۇ n/كۈلمەك nv/بولۇۋىدى v/شەيخ/نول/كۆزىنى n/قىسىپ dv/ئۆك/اقىشىنى/ئاتتى v/.
3. ئۇ/ئەمدى y/ئاتىسىغا n/گىشىپ dv/سالام/قىلىشتىنمۇ nv/قالدى v/.
4. بۇ/ئىك/ئالاقانات n/قاپ/يۈرەك n/يىگىتتىن/قانداق rs/ئىنتىقام n/ئېلىش nv/ئۈستىدە/ئويلىدى v/.

2.2 سۆز تۈركۈمى سىتاسىتىكا ئۇچۇرلىرىغا ئېرىشىش مودىلى

سۆز تۈركۈمى ئېھتىماللىق جەدۋىلى ۋە سۆز تۈركۈملىرىنىڭ نىسپەت جەدۋىلى كىيىنكى باسقۇچلاردىكى سۆز تۈركۈمىنى بەلگىلىگۈچ مودىلىدا ئىشلىتىلىدۇ. بۇ بۆلەكتە ئۇلارنىڭ ھىساپلاش ئۇسۇلى قىسقىچە تونۇشتۇرىلىدۇ.

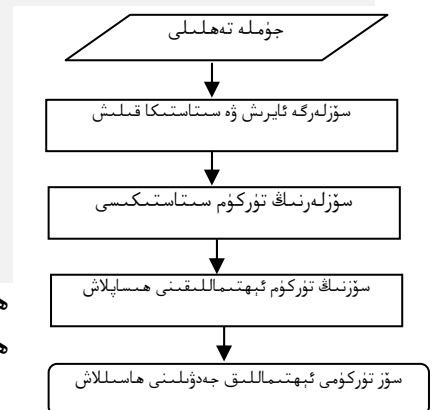
2.2.1 سۆز تۈركۈمى ئېھتىماللىق جەدۋىلى

VisualStudio2008 مۇھىتىدا C# پروگرامماتىلى ئارقىلىق سۆز تۈركۈمى ئېھتىماللىق جەدۋىلىگە ئېرىشىمىز، بۇ جەدۋەلنى ھاسىل قىلىش جەريانىدا يېزىلغان يادرولۇق كود بۆلىكى تۆۋەندىكىچە:

```
foreach (Sentence sen in Sentences)
{
    foreach (Word wrd in sen.Words)
    {
        if (TotalWords.ContainsKey(wrd.Content)) // سۆزنىڭ كۆرۈلۈش قېتىم سانىنى ھېساپلاش
        {
            TotalWords[wrds.Content]++;
        }
        else
        {
            TotalWords.Add(wrd.Content, 1);
        }
        if (pools.Keys.Contains(wrd.Content + wrd.Turkum)) //
        {
            pools[wrds.Content + wrd.Turkum].Count++;
        }
        else
        {
            pools.Add(wrd.Content + wrd.Turkum, new WordPool(wrd, 1));
        }
    }
}
// يېڭىدىن كۆرۈلگەن سۆز ۋە ئۇنىڭ سۆز تۈركۈمىنى خاتىرىلەش
nCount++;
```

ھاسىل قىلىنغان سۆز تۈركۈمى ئېھتىماللىق جەدۋىلىنىڭ مەزمۇنى ۋە ھىساپلاش ئۇسۇلى 2-رەسىمدە كۆرسىتىلگەندەك:

```
ProbWord[] = {{"ikki", "m", 0.0020703933747422},
{"d0rE", "q", 0.5},
{"mExrEptin", "n", 0.5},
{"keyin", "dv", 0.00429184549356223},
```



2-رەسىم سۆز تۈركۈمى ئېھتىماللىق جەدۋىلى دىئاگراممىسى

.....

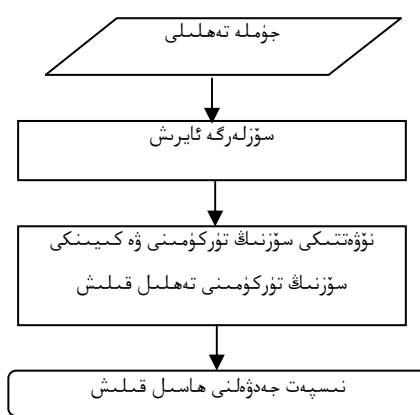
2.2.1 سۆز تۈركۈمى نىسپەت جەدۋىلى

VisualStudio2008 مۇھىتىدا تۇرۇپ C# تىلى ئارقىلىق پروگرامما تۈزۈش ئارقىلىق سۆز تۈركۈملىرى بەلگىلەنگەن تېكىستىن سۆز تۈركۈمى نىسپەت جەدۋىلىگە ئېرىشىمىز، بۇ جەدۋەلنى ھاسىل قىلىش جەريانىدا يېزىلغان يادرولۇق كود بۆلىكى تۆۋەندىكىچە:

```
for (i = 0; i < 38; i++)
{sb.Append("{");
sbHtm.AppendLine("<tr>");
for (j = 0; j < 38; j++)
{n4[i, 38] = n4[i, 38] + n4[i, j]; } // ھەر قايسى سۆز تۈركۈملىرىنىڭ كۆرۈش قېتىم سانىنى ھېساپلاش
nTotal = nTotal + n4[i, 38];
for (j = 0; j < 39; j++)
{sb.Append(n4[i, j]);
sbHtm.AppendLine("<td>" + n4[i, j] + "</td>");
if (j < 38)
{sb.Append(", ");}}
```

ھاسىل قىلىنغان سۆز تۈركۈمى نىسپى جەدۋىلىنىڭ مەزمۇنى ۋە ھىساپلاش ئۇسۇلى 3-رەسىمدە كۆرسىتىلگەندەك:

```
Int BiGramProbs[37][38]={\{0,0,0,0,0,0,0,0,0,0,0,0,  
0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,  
0,0,0,0,0,0\}  
\{0,987,16,173,39,56,75,43,29,4,161,8,1,5,8,5016,92,  
123,3,56,0,1,53,5,105,10,68,24,5,8,18,32,1,881,578,  
175,427,44,32,0,1,0,9363\}}
```



3-رەسىم سۆز تۈركۈملەرنىڭ نىسپەت جەدۋىلى دىئاگراممىسى

2.3 سۆز تۈركۈمىنى ئاپتوماتىك بەلگىلىگۈچى مودىلى

ئاپتوماتىك بەلگىلەشنى ئەمەلگە ئاشۇردىغان ئالگورىتملارنى ئىككى تۈرگە ئايرىشقا بولىدۇ؛ ئۇلار: مەلۇم قائىدىلەر ئاساسىدىكى ئۇسۇل (rule-based tagger) ۋە رەندىم ئايرىش ئۇسۇلى (stochastic tagger). رەندىم ئايرىش ئۇسۇلىدا ئادەتتە مەشقىلەندۈرۈش ئامبىرى ئارقىلىق بەلگىلەش ئوبيكتى بولغان ھەر بىر سۆزنىڭ جۈملە مۇھىتىدىكى تەۋەلىك تىپ ئېھتىماللىقى ھىساپلاپ چىقىلىدۇ. مەسلەن: HMM بەلگىلەش سىستېمىسى [2].

بۇ ماقالىدە ئىنتوناتسىيە قەۋەتلىرىنى ئاپتوماتىك ئايرىش ئېھتىياجىگە ئاساسەن HMM بەلگىلەش مۇدىلىدىكى ئىككى ئىلمىنلىق ئۇسۇل قوللۇنۇلۇپ ئۇيغۇرتىلى سۆز تۈركۈمىنى ئاپتوماتىك بەلگىلەش سېستىمىسىنىڭ ئاساسى گەۋدىسى ھاسىللاندى. HMM بەلگىلەش مۇدىلىدىكى ئىككى ئىلمىنلىق ئۇسۇلنىڭ ھىساپلاش جەريانى تۆۋەندىكىچە:

(1) ". بەلگىسى بولسا جۈملە ئاخىرلاشقانلىقىنىڭ ئىپادىسى بولۇپ،يۇنىڭدىن پايدىلىنىپ بىر جۈملە تاللاپ ئېلىۋېلىندۇ.

(2) مەشقلەندۈرۈلدىغان جۈملىلەردىن ھەر قايسى جۈملىلەر ئايرىم-ئايرىم تاللاپ ئېلىنىپ بۇ جۈملىلەرنى تەشكىللەنگۈچى سۆزلەرنىڭ سۆز تۈركۈمى نىسبىتى ھىساپلاپ چىقىلىدۇ.

$$P(w_i | c_i) = \frac{\text{训练语料中 } w_i \text{ 的词性被标记为 } c_i \text{ 的次数}}{\text{训练语料中 } c_i \text{ 出现的总次数}} \quad (1)$$

فورمىلاسى :

سۈرەت: مەشقىلەندۈرۈش توپىدا w_i سۆزىنىڭ سۆز تۈركۈمى c_i دەپ بەلگىلىنىش قېتىم سانى

مەخرەج: مەشقىلەندۈرۈش توپىدا c_i نىڭ جەمئىي كۆرۈلۈش قېتىم سانى

(3) نۆۋەتتىكى سۆزنىڭ ئامباردا كۆرۈلگەن بارلىق سۆز تۈركۈملىرى تارتىپ چىقىرىلىپ تۈگۈن نوقتتا ھاسىل قىلىنىدۇ، ئاندىن ھەر بىر سۆز تۈركۈمىنىڭ ئۆزىنىڭ ئالدىدىكى سۆزنىڭ ھەر بىر سۆز تۈركۈمى بىلەن بولغان بارلىق گۇرۇپپىلىنىش ئىھتىماللىقى ھىساپلىنىدۇ، ھەمدە بۇ ئارقىلىق بىز يۇقاردا تونۇشتۇرۇپ ئۆتكەن سۆز تۈركۈملىرىنىڭ نىسبىي جەدۋىلى ھاسىللىنىدۇ.

$$P(c_i | c_{i-1}) = \frac{\text{تەلىم تىلىمدا } c_i \text{ نىڭ } c_{i-1} \text{ نىڭ ئالدىدىكى سانى}}{\text{تەلىم تىلىمدا } c_{i-1} \text{ نىڭ ئالدىدىكى سانى}} \quad (2)$$

فورمىلاسى :

سۈرەت: مەشقىلەندۈرۈش توپىدا c_i نىڭ c_{i-1} دىن كىيىن كېلىدىغان ئەھۋالنىڭ كۆرۈلۈش قېتىم سانى

مەخرەج: مەشقىلەندۈرۈش توپىدا c_{i-1} نىڭ جەمئىي كۆرۈلۈش قېتىم سانى

(4) سۆزنىڭ تۈركۈمىنى بېكىتىشتە، جۈملىدىكى سۆزلەرنىڭ تەرتىپى بويىچە باشلىنىش نوقتىسىغا باشلىنىش تۈگۈن نوقتىسىنى قويۇۋالىمىز. بىرىنچى سۆزنىڭ بارلىق ئىھتىماللىق سۆز تۈركۈملىرىنى تارتىپ چىقىرىۋېلىپ ھەر بىر سۆز تۈركۈمىگە قاراپ تۈگۈن نوقتىسى قويۇۋالىمىز، بىرىنچى سۆزنىڭ ھەر بىر سۆز تۈركۈمى تۈگۈن نوقتىسىنى باشلىنىش تۈگۈن نوقتىسىغا يۈزلەندۈرۈپ كىيىنكى سۆزنى چىقىرىۋالىمىز.

(5) HMM ئاساسىدىكى بەلگىلەش ئالگورىتمىدا سۆزنىڭ ئىككى ئىلمىنىلىق ئىھتىماللىقى بىلەن گۇرۇپپىلىنىش ئىھتىماللىقىنىڭ كۆپەيتىمىسى ئەڭ چوڭ بولغان بەلگىلەر ئارقىمۇ-ئارقىلىقى ئەڭ ئاخىرقى يەكۈن قىلىنىدۇ^[2]. يەنى $P(w_i | c_i) P(c_i | c_{i-1})$ ھىساپلىنىپ، نەتىجىسى ئەڭ چوڭ بولغان ئىككى ئىلمىنىلىق سۆز تۈركۈمى گۇرۇپپىسى تېپىپ چىقىلىدۇ، ھەمدە نۆۋەتتىكى سۆزنىڭ سۆز تۈركۈمى تۈگۈن نوقتىسىنىڭ ئالدىدىكى ئەڭ ئۈنەملىك تۈگۈن نوقتىسى بىلەن بىرلىكتە، بۇ ئىككى ئىلمىنىلىق گۇرۇپپىنىڭ كەينىدىكى سۆز تۈركۈمى تۈگۈن نوقتىلىرىنى كۆرسىتىدۇ.

(6) ئەگەر بىر تەرەپ قىلىنمىغان سۆز بولسا بەشىنچى باسقۇچ بويىچە داۋاملىق ھىساپلاش ۋە بەلگىلەش ئېلىپ بېرىلىدۇ، ئەگەر ئۇنداق سۆزلەر بولمىسا كىيىنكى باسقۇچ ئورۇندىلىدۇ.

(7) بارلىق سۆزلەر بىر تەرەپ قىلىنىپ بولۇنغاندىن كىيىن جۈملىنىڭ ئەڭ كەينىدىكى بىر سۆزدىن باشلاپ ھەر بىر سۆزنىڭ ئۆزىگە خاس بولغان بارلىق سۆز تۈركۈملىرى ئىچىدىكى ئەڭ مۇۋاپىق بولغان سۆز تۈركۈمى لىنىيەسى بويىچە بۇ جۈملىدىكى ھەر بىر سۆزنىڭ سۆز تۈركۈمى ئارقىمۇ-ئارقىلىقى تېپىپ چىقىلىدۇ.

VisualStudio2008 مۇھىتىدا تۈرۈپ C# تىلى ئارقىلىق پروگرامما تۈزۈش ئارقىلىق بۇ بۆلەك ئاچقۇچلۇق ئالگورىتم بۆلىكىدىن ئارزۇ قىلغىنىمىزدىكىدەك نەتىجىگە ئېرىشىمىز، بۇ جەرياندا يېزىلغان يادرولۇق كود بۆلىكى تۆۋەندىكىچە:

```
for (int nIndex = strWords.Count() - 1; nIndex > -1; nIndex--)
{
    if (nIndex == strWords.Count() - 1)
    {
        if (lstWords.Where(w => w.Content == strWords[nIndex]).Count() > 0)
        {
            var WordResult = lstWords.Where(w => w.Content == strWords[nIndex]).GroupBy(w => w.Turkum).OrderByDescending(w => w.Count()).First();
        }
    }
}
```

```

straResultWords[nIndex] = straWords[nIndex] + "/" + WordResult.Key;}
else
{straResultWords[nIndex] = straWords[nIndex] + "/n";}}
else
{List<NisbetJedwili> lstCurrentCollection = new List<NisbetJedwili>();
var CurrentCollection = lstNisbet.Where(n => n.Content == straWords[nIndex]);
foreach (var c in CurrentCollection)
{lstCurrentCollection.Add(c);}
List<NisbetJedwili> lstNextCollection = new List<NisbetJedwili>();
var NextCollection = lstNisbet.Where(n => n.Content == straWords[nIndex + 1]);
foreach (var n in NextCollection)
{lstNextCollection.Add(n);}
straResultWords[nIndex] = straWords[nIndex] + getTurkum(lstCurrentCollection,
lstNextCollection);}

```

3. تەجربە مۇھىتى ۋە ئانالىزى

3.1 سانلىق مەلۇماتلارنى تەييارلاش

تەجربە سانلىق مەلۇماتلىرىنى ھازىرلاش جەريانىدا مەشىقلەندۈرۈش سانلىق مەلۇماتلىرى ۋە سىناق سانلىق مەلۇماتلىرى ھازىرلاندى. بۇنىڭدا گېزىت-ژورنال، ئەدەبىي ئەسەرلەر قاتارلىق مىدىئا ۋاسىتىلىرىدىن تەكست ماتېرىيالى يىغىۋېلىنىپ ئاندىن تەكست تاللاشتىكى بىر قانچە باسقۇچلار ئارقىلىق 10610 دانە جۈملە تاللىۋېلىنىپ بۇنىڭ ئىچىدىكى 10000 جۈملە مەشىقلەندۈرۈشكە قاتناشتۇرۇلدى. بۇ جۈمىلەلەرنىڭ كۈندۈلۈك تۇرمۇشتىكى قوللىنىلىش نىسپىتى بىر قەدەر يۇقىرى. سىناق جۈمىلىرى توپلام ئىچىدىكى سىناق جۈمىلىسى ۋە توپلام سىرتىدىكى سىناق جۈمىلىسى دەپ ئىككى تۈرگە ئايرىلغان بولۇپ ئايرىم-ئايرىم ھالدا 500 جۈملىدىن تەركىپ تاپتى. بۇ سانلىق مەلۇماتلار ئۆزى تەۋە بولغان توپلامدىن ئىختىيارى تاللىۋېلىندى.

سۆز تۈركۈمىنى بەلگىلەش نېگىزلىك ئېيتقاندا كۆپ خىللىق مەسلىسى (disambiguation) نى ھەل قىلىش بولۇپ ھازىرقى ئامبار (Alfa-ئامبار) دا جۈملە مۇھىتىدىكى ھەر بىر سۆزگە بىردىن سۆز تۈركۈمى بەلگىلەندى، جەمئى 123678 سۆزدىن تەركىپ تاپقان بولۇپ بۇنىڭ ئىچىدىكى %50.39 سۆز تۈركۈم جەھەتتە ھەرخىللىققا ئىگە. 2-جەدۋەلدا Alfa-ئامباردىكى سۆزلەرنىڭ ھەرخىللىق ئەھۋالى ۋە مۇناسىۋەتلىك سانلىق مەلۇماتلار كۆرسىتىلدى.

2-جەدۋەل Alfa-ئامباردىكى سۆزلەرنىڭ ھەرخىللىق ئەھۋالى ۋە مۇناسىۋەتلىك سانلىق مەلۇماتلار

تەكرارلانمىغان سۆزنىڭ ئومۇمى سانى	ھەرخىللىق ئەھۋالى
25412	بىرلا خىللىق (بىرلا تۈركۈمگە تەۋە بولغىنى)
3183	ھەرخىللىق (2-7 خىلغىچە تۈركۈمى بولغىنى)
2497	ئىككى تۈركۈمگە تەۋە بولغىنى
468	ئۈچ تۈركۈمگە تەۋە بولغىنى
165	تۆت تۈركۈمگە تەۋە بولغىنى
52	بەش تۈركۈمگە تەۋە بولغىنى

بوغۇمنىڭ سىناق سانلىق مەلۇماتلارغا بولغان تەسىرى كۆزدە تۇتۇلۇپ تەجرىبە جەريانىدا بوغۇم ئامىلىمۇ نەزەردىن ساقىت قىلىنمىدى، يەنى ئالدى بىلەن Alfa ئامبىرىغا قارىتا بوغۇم سىتاتىستىكىسى ئېلىپ بېرىلدى. مەزكۇر ئامباردا بىر سۆز كۆپ بولغاندىمۇ 10 دانە بوغۇمدىن تەركىپ تاپقان بولۇپ، ئەڭ ئاز بولغاندا بىر بوغۇمدىن تەركىپ تاپقان. توققۇز ۋە ئون دانە بوغۇمدىن تەركىپ تاپقان سۆزدىن پەقەت بىر تالدىنلا بار، بوغۇم سىتاتىستىكا ئۇسۇلى ئامباردىكى سانلىق مەلۇماتلارنىڭ خاتالىقلارنى تۈزۈتۈپ ساپلىقنى ئاشۇرۇشتا مۇھىم رول ئوينىدى.

3.2 تەجرىبە سۇپىسى ۋە نەتىجە ئانالىزى

تەجرىبە جەريانىدا مىكرو-سوفىت شىركىتى ئوتتۇرغا قويغان ئەڭ يېڭى پروگرامما تۈزۈش تىلى بولغان C#، تېخنىكىلىق ئېچىش تىلى سۈپىتىدە قوللىنىلغان بولۇپ، ئۇ ئىشلىتىشكە ئەپلىك بولۇش، پۈتۈنلەي ئابونتقا يۈزلەنگەن بولۇش ۋە تۈر جەھەتتە بىخەتەر بولۇشتەك ئالاھىدىلىكلەرگە. بۇ ئالاھىدىلىكلەر ئۇنى كىيىنكى بىر دەۋىردىكى تەقسىمات تىپى قوللىنىشچان پروگراممىسىنىڭ ئاساسىي ئېقىم ئىجادىيەت تىلىغا ئايلاندۇرغان [7]. Net Framework 3.5. بولسا Linq تېخنىكىسىغا ئىگە بولۇپ Unicode بىر تەرەپ قىلىشقا ناھايىتى ماس كېلىدۇ. پروگرامما ئېچىش تىلى بولسا visual studio 2008 .steam system.

ئېنتوناسىيە قاتلاملىرىنى ئاپتوماتىك بەلگىلەش تەتقىقاتى جەريانىدا ئەتراپتىكى سۆزلەرنىڭ سۆز تۈركۈم ئالاھىدىلىكى، نوۋەتتىكى سۆز ۋە كىيىنكى سۆزدىكى جۈملىنىڭ ئۇزۇنلۇق ئالاھىدىلىكى، جۈملە بېشىدىن باشلاش ياكى جۈملە ئاخىرىدىن باشلاشتىكى ئالدىنقى چىگرا ئارلىق تىپى، ھەمدە ئەتراپتىكى سۆزدە ئۇرغۇنىڭ بولۇش بولماسلىقى قاتارلىق ئالاھىدىلىكلەر بولسا نوقتىلىق ئويلىشىلىدىغان ئىلمىنلار بولۇپ ھىساپلىنىدۇ [2]. سىناق باسقۇچىدا تەييارلىنىدىغان ئىككى خىل سىناق توپلىمىنىڭ بىرى مەشىقلەندۈرۈش توپلىمىدىن ئېلىنغان ئىچكى سىناق سانلىق مەلۇماتى، ئۇلار توپلامدىن ئىختىيارى تاللاپ ئېلىنغان، يەنە بىرسى بولسا مەشىقلەندۈرۈش توپلىمىنىڭ سىرتىدىن تاللىۋېلىنغان سىناق سانلىق مەلۇماتلاردۇر. سېستىمىنىڭ ئىجرا بولۇش كۆرنىشى تۆۋەندىكىچە:

4-رەسىم ئىجرا بولۇش

سىناقنىڭ توغرىلىق دەرىجىسىنى كۆزىتىشتە توغرىلىق نىسپىتى ۋە چاقىرىش نىسپىتى ھىساپلاندى. توغرىلۇق نىسپىتى = سۆز تۈركۈمى توغرا بەلگىلەنگەن سۆزنىڭ ئۇمۇمى سانى / سىناق تېكىست تەركىبىدىكى سۆز تۈركۈملىرىنىڭ ئۇمۇمى سانى

ھەر قايسى سۆز تۈركۈمىنىڭ چاقىرىش نىسپىتى=شۇ تۈركۈمگە تەۋە بولغان سۆزلەرنىڭ توغرا ئايرىغان ئۇمۇمى سانى\سنىق توپلامدىكى مەزكۇر سۆز تۈركۈمىنىڭ ئۇمۇمى سانى
چاقىرىش نىسپىتىنىڭ ئورتالانمىسى=ھەر قايسى سۆز تۈركۈملىرىنىڭ خاس چاقىرىش نىسپەتلىرىنىڭ يىغىندىسى\سۆز تۈركۈمىنىڭ سانى

3 -رەسىم سنىق نەتىجىسى

ئۇسۇل باھالاش ئۆلچىمى	مەشىقلەندۈرۈش توپلىمىنىڭ ئىچىدىن تاللىۋېلىنغان سنىق سانلىق مەلۇماتلىرى 500 جۈملە	مەشىقلەندۈرۈش توپلىمىنىڭ سىرتىدىن تاللىۋېلىنغان سنىق سانلىق مەلۇماتلىرى 500 جۈملە
توغرۇلۇق نىسپىتى	86.92%	71.99%
چاقىرىش نىسپىتى	71.66%	62.42%

سۆز تۈركۈمى 41 خىلدىن 37 خىلغا ئازلىغاندا توپلام ئىچى سىنىقىنىڭ توغرۇلۇق نىسپىتى 86.69% دىن 86.92% غا، چاقىرىش نىسپىتى 65.81% دىن 71.66% غا يۇقىرى كۆتۈرۈلدى. توپلام سىرتى سىنىقىنىڭ توغرۇلۇق نىسپىتى بولسا 72.12% دىن 71.99% غا تۆۋەنلەپ، چاقىرىش نىسپىتى 54.61% دىن 62.42% غا كۆتۈرۈلدى. گەرچە 37 خىل سۆز تۈركۈمى ئاساسىدىكى توپلام سىرتى سىنىقىنىڭ توغرۇلۇق دەرىجىسى سەل تۆۋەنلىگەن بولسىمۇ ئەمما باشقا قىممەتلەردە كۆرۈنەرلىك دەرىجىدە يۇقىرى ئۆرلەش بولغان. تەجرىبە جەريانىدا ئېرىشلىگەن نەتىجە قوللىنىش تەلپىگە ئۇيغۇن كېلىدۇ. شۇڭلاشقا 37 خىل سۆز تۈركۈمى ئاساسىدىكى ئىككى ئىلمىنىڭ بولغان ئاپتوماتىك سۆز تۈركۈمىنى بەلگىلەش ئالگورىتمى، ئېنتوناسىيە قەۋەتلىرىنى ئاپتوماتىك ئايرىش ئالگورىتمىنىڭ ئاساسى سۈپىتىدە قوللىنىلسا بولىدۇ.

4. ئاخىرلاشتۇرۇش سۆزى

بۇ ماقالىدە ئېنتوناسىيە چىگرىسىنى ئاپتوماتىك ئايرىش تېخنىكىسىغا يۈزلەندۈرۈلگەن ئۇيغۇر تىلىدىكى سۆز تۈركۈملىرىنى ئاپتوماتىك بەلگىلەش تېخنىكىسىغا قاراپ تەتقىقات ئېلىپ بېرىلدى. بۇ قېتىملىق ئىزدىنىش ئارقىلىق، HMM مودىلى ئىككى ئىلمىنىڭ ئۇسۇلىنىڭ، ئۇيغۇر تىلىدىكى ئېنتوناسىيە چىگرىسىنى ئاپتوماتىك ئايرىش جەريانىدا ئىشلىتىلىدىغان سۆز تۈركۈمىنى ئاپتوماتىك بەلگىلەش تەلپىگە ئۇيغۇن كېلىدىغانلىقى ئىسپاتلاندى. گەرچە بىر قىسىم خاتالىقلار مەۋجۇت بولسىمۇ ئەمما ئاساسى مەخسەت بولغان ئېنتوناسىيە چىگرىسىنى ئاپتوماتىك ئايرىش تەلپىگە ئۇيغۇن ۋە ئېرىشلىگەن نەتىجە ئېھتىياجىنى قاندۇرالايدۇ. ئەمما ھىساپلاش ئۇسۇلىغا نىسپەتەن يەنىمۇ بىر قەدەر ئىلگىرلىگەن ئاساستا ئىسلاھات ئېلىپ بېرىش كېرەك. ئۇندىن باشقا يەنە مەشىقلەندۈرۈلدىغان سانلىق مەلۇماتلار سانىنى ئاشۇرۇش، سۆز تۈركۈملىرىنىڭ سانى قاتارلىق مەسىللەردە يەنىمۇ بىر قەدەر ئىلگىرلىگەن ھالدا ئانالىز ئېلىپ بېرىش كېرەك.

参考文献

- [1] 童方. 语音合成系统的关键技术与应用实例[J]. 科海拾贝, 2000. (2):19-21.
- [2] Daniel Jurafsky, James H. Martin. Speech and Language Processing[M]. Beijing: Publishing house of Electronics Industry, 2005.
- [3] 艾斯卡尔·肉孜. 基于HMM的维吾尔语音合成系统的研究与实现[D]. 新疆乌鲁木齐: 新疆大学, 2008. 7.
- [4] 洪梅, 吐尔根·伊不拉音. 维吾尔语的词性标注对校初谈[J]. 微型电脑引用, 2006. (22):12.
- [5] 田斌, 易克初. 一种可扩展的汉语统计二元文法语言模型[J]. 信号处理, 2002. (3):184-187.
- [6] Francis, W.N. and Kučera, H. (1982). Frequency Analysis of English Usage. Houghton, Boston.
- [7] 康博. c#程序员参考手册[M]. 北京:清华大学出版社, 2002.
- [8] 赵晟, 陶建华, 蔡莲红. 基于规则学习的韵律结构预测[J]. 中文信息学报, 2002, 16(5):30-37.

- [9] 吴义坚. 基于隐马尔科夫模型的语音合成技术研究[D]. 中国科学技术大学, 2006. 4.
- [10] 杨行峻, 迟惠生. 语音信号数字处理[M]. 电子工业出版社, 1995.
- [11] 荀恩东, 钱揖丽等. 应用二叉树剪枝识别韵律短语边界[J]. 中文信息学报, 2005. 20(3):24-28.
- [12] 陈永彬, 王仁华. 语音信号处理[M]. 中国科技大学出版社, 1990.
- [13] 吴志勇, 蔡莲红. 语音合成中的韵律关联模型[J]. 中文信息学报, 2003, 18(2):44-50.