

ՀԱՇՎՈՂԱԿԱՆ ՏԵԽՆԻԿԱՑԻ ԿԻՐԱՌՈՒՄԸ ՏԵՔՍԻ
ՈՐՈՇԱԿԻԱՑՄԱՆ ՀԱՐՑՈՒՄ

Ինֆորմացիայի տեսության մեջ կապակցված տեքստը դիտվում է ինֆորմացիայի էրգոդիկ աղբյուրու հսկ վերջինիս հատուկ են վիճակագրական որոշակի օրինաչափություններ, որոնք ենթարկվում են բաշխվածության օրենքի և հաղորդագրության երկարության ալիլացման միջոցով կարելի է հասնել կայուն ու ճշգրիտ գնահատականների: Այդ նշանակում է, որ որևէ լեզվով գրված ամեն մի տեքստի համար կարող ենք ստանալ բնութագրիչների համախումբ, որի օգնությամբ մի տեքստը տարբերվում է մեկ այլ տեքստից:

Դժվար է պնդել, թե ամեն մի տեքստ պատկանում է ինֆորմացիայի իդեալական էրգոդիկ աղբյուրների թվին, մանավանդ, երբ խոսքը վերաբերում է գրական հուշարձաններին: Ցանկացած տեքստում առկա են լեզվական ընդհանուր ու հեղինակային անհատականությունը, դարաշրջանային համաժամանակությունն ու տարժամանակությունը, հնարանությունն ու նորաբանությունը, և եթե այս ամենին ավելացնենք այն, որ գրական հուշարձանները մեզ հասել են ոչ թե բնագրային տեսքով, այլ տարբեր ժամանակներում և տարբեր կրթօջախներում արված կրկնօրինակների, այսպես կոչված ցուցակների (կամ ընդօրինակությունների) տեսքով, որոնք իրենց վրա կրելով դարաշրջանային, տեղային աղեցությունները, պարունակում են օրինաչափությունների միօրինակների և օրինակների բնույթի աղավաղումների: Վերջինները կարող են հասնել այն ասակին անհամար կամ առնվի տվյալ ընդորինակությունը տվյալ հեղինակին պատկանելու հարցը: Այդ աղավաղումների պատճառով առանձնապես դժվարանում է տեքստային բնութագրիչների որոշակի և դրանց գնահատումը:

Ընդհանրապես, տեքստաշափական աշխատանքներին հատուկ են տեքստային լայն գանգվածների ներառմամբ գործողությունների միօրինակություն, երբ անհրաժեշտ է հաշվի առնել հարյուրավոր փաստեր, որը հաճախ գերազանցում է մարդու հնարավորությունները: Մյուս կողմից, ժամանակակից հաշվողական տեխնիկան թույլ է տալիս ընդգրկել ավելի

լայն տեքստային գանգվածներ, քննության առնել ամեն մի փաստ, առաջ քաշել խնդիրներ, որոց վճռումը ավանդական եղանակով գործնականում հնարավոր չէ:

Այս առումով ներկա հոդվածում դրվում է երկու խնդիր կապված Մ. Խորենացու «Հայոց պատմություն» ստեղծագործության քննության հետ և նշվում այդ խնդիրների լուծման եղանակները ժամանակակից էլեկտրոնային հաշվողական մեքենաների օգնությամբ:

Ժամանակակից բանասիրության մեջ գրական ստեղծագործությունների որոշակիացման խնդիրը քննվում է երեք ուղղություններով՝ փաստավեճրագրական գաղափարաթեմատիկ և լեզվառնական: Այս կապակցությամբ Վ. Վ. Վինոգրադովը գրում է. «Ոճական վերլուծության մարքսիստական մեթոդաբանությունը, մեր կարծիքով, պետք է բխի գրական ստեղծագործության իմացությունից՝ ինչպես գաղափարական, թեմատիկ, կոմպոզիցիոն և լեզվական (ընդգծումը մերն է) դիալեկտիկորեն կապված տարրերի կառուցվածքային միասնությունից» (Վ. Վ. Վինոգրադով, էջ 197):

Հայտնի է, որ ինչպես տեքստի, այնպես էլ նրա կաղմի մեջ մտնող ամեն մի նախադասության տարրերը, ոչ թե մեկուսացված, տարանջատ այլ համակարգայնության հարաբերություններով կապված լեզվական միավորներ են, որոնք հանդիս են գալիս որոշակի օրենքների ու կանոնների համապատասխան: «Ամեն մի գործողություն,— գրում է Զ. Շմիդտը, — լինի դա լեզվական, թե ոչ լեզվական, նկատի ունի որոշակի նորմա և տեխնիկա, որոնք առանց էական շեղումների ընդունվում ու գործառվում են հասարակության բոլոր անդամների կողմից» (Շմիդտ, էջ 100):

Լայն ըմբռնմամբ նորմայի ու տեխնիկայի համախումբը կարելի է անվանել քերականություն, որը, փաստորեն, տրված ելակետային միավորների հետ նախօրոք սահմանված գործողությունների հաշորդակարգ է: Եթե ոչ բնական իրողություններում օրյեկտների միջև գոյություն ունեցող հարաբերությունները կարելի է նկարագրել այսպես կոչված սոցիալական քերականության օգնությամբ, ապա լեզվական իրողություններում լեզվական միավորների միջև հարաբերությունները նկարագրվում են լեզվի քերականության՝ շարահյուսության և ձևարանության օգնությամբ:

Այսպիսով, տեքստի բառակազմի ընտրությունը և նրանում բառերի հաջորդման կարգը հիմնականում որոշվում է լեզվական նորմայով: Սակայն լեզվական նորման քարացած հրահանգների հավաքածու չէ: Լեզվական նորմայի կիրառության վրա էական ազդում է հեղինակային անհատականությունը, որով և բնորոշվում է ոճը: Հետևաբար, տեքստի հե-

ղինակային ոճի որոշումը հանգում է լեզվական նորմայի կիրառման առանձնահատկությունների բացահայտմանը:

Տեքստի վիճակագրությունը ցուց է տալիս, որ Հեղինակային ոճի որոշման հարցում տեքստի առանձին միավորների արձանագրումը խիստ կարևոր է, սակայն հուսալի արդյունքներ չեն ապահովում: Արդյունքներն ավելի հավաստի ու ծանրակշիռ են դառնում, երբ իրեւ հետազոտության օբյեկտը ընտրվում են ոչ միայն տեքստի առանձին միավորները, այլև նրանց տարրերը զուգորդումները: Հարկ է նշել, որ քննության համար տեքստից կարող են ընտրվել ինչպես բառային միավորները, այնպես էլ դրանք բնութագրող քերականական կարգերը:

Սահմանում Տեքստում լեզվական A միավորի (կամ նրա բնութագրող քերականական կարգի) հանդես գալը անվանենք լեզվական միավորի (կամ քերականական կարգի) հանդիպում և նշանակենք A(п)-ով, որտեղ ո-ը A միավորի հանդես գալու թիվն է:

Լեզվական A և B միավորների (կամ նրանց քերականական կարգերի) հարագիր ձևով տեքստում հանդես գալը անվանենք A և B միավորների զույգ հանդիպում և նշանակենք A ո B, որտեղ ո-ը այնպիսի հարագիրության հանդիպումների թիվն է, երբ A-ն գրավում է ձախադադիրը:

Տեքստում լեզվական միավորների հանդիպման քննությունը հանդեցնում է միաշափ աղյուսակի, որտեղ տրված են լեզվական միավորները և հանդիպման ո-հաճախությունները. այդ միավորների զույգ հանդիպումների քննությունը՝ երկշափ աղյուսակի՝ քառակուսի մատրիցի, որի տողերն ու սյուները փոխհամապատասխանության մեջ են դրվում A և B լեզվական միավորների հետ, իսկ մատրիցի յուրաքանչյուր տարրը ներկայացնում է զույգ հանդիպման թիվը կամ հաճախությունը:

Տեքստի կարեւոր բնութագիրներից է բառաձևերը բնորոշող քերականական կարգերի կամ նրանց համախմբերի փոխդասավորությունը: Այն որոշելու համար տեքստի յուրաքանչյուր բառը փոխարինվում է նրան համադրվող քերականական կարգերի համախմբով, գոյական՝ անուն, հոլով, թիվ, հոլովում, բայց եղանակ, ժամանակ, թիվ, դեմք և այլն:

Մ. Խորենացու «Հայոց պատմության» քննության համար ընտրել ենք քերականական կարգերի 360 համախմումք:

Գոյական գն <անուն, հոլով, թիվ, հոլովում>,

անուն <հատուկ (հտ), հասարակ հս>,

հոլով <ուղղական (ուղղ), սեռական (սռ), արական (ար),

հայցական (հյ), բացառական (բց), գործիական (գծ)>,

թիվ <եզակի (եթ), հոգնակի (հթ)>,

նոլովում <ա, ե, ի, ո, ու, ի-ա (իա), ո-ա, (ոա)>,

ածական ան <ածականի տեսակ, հոլով>,
ածականի տեսակ <դրական (դր), համեմատական (հմ), գերա-
դրական (գր), հարաբերական (հր)>,
թվական թե <թվականի տեսակ, հոլով, թիվ>,
թվականի տեսակ <քանակական (քն), դասական(դս), անձներա-
կան (աձ)>,
դերանուն դն <դերանվան տեսակ, հոլով>,
դերանվան տեսակ <անձնական (ան), փոխադարձ (փձ), ցուցա-
կան (ցց), ստացական (սց), հարցական (հց) հա-
րաբերական (հր), անորոշ (աշ), որոշյալ(ոշ)>,
բայ բյ <եղանակ, ժամանակ, թիվ, դեմք>,
եղանակ <սահմանական (սհ), ստորագասական (ստ), հրամայա-
կան (հր), սահմանական եղանակի ժամանակ ներկա (նկ),
անցյալ կատարյալ (ակ)>,
ստորադասական եղանակի ժամանակ <առաջին ապառնի (առ),
երկրորդ, ապառնի (եր)>,
հրամայական եղանակի ժամանակ <բուն (բն), հորդորական (հն),
արգելական (ար)>,
դեմք <առաջին (ադ), երկրորդ (բդ), երրորդ (գդ),
դերբայ դր <տեսակ, հոլով>,
դերբայի տեսակ <անորոշ (աշ), անցյալ (աց), ապառնի (ապ)
ենթադրական (են)>,
շաղկապ շղ <համադասական (հդ), ստորագասական (սդ)>,
մակրայ մր <տեղի (տղ), ժամանակի (ժմ), ձեփ (ձլ), շափի (շփ),
աստիճանի (աս)>,
Ժխտական մասնիկ մս,
նախդիր նդ <հոլով>,
նախադրույուն նիւ <հոլով>,
կետադրույթյուն <միջակետ (մտ), ստորակետ (ստ), բութ (բթ),
վերջակետ (վշ)>;
Քերականական կարգերի ընտրված համակարգում գոյականին հա-
մադրվում է 168 համախումբ (1-168), ածականին՝ 24 (168-12), թվա-
կանին՝ 36 (298-228), դերանվանը՝ 48 (229-276), բային՝ 36 (277-312),
դերբային՝ 24 (213-236), շաղկապին՝ 2 (237-238), մակրային՝ 5 (338-
343), ժխտական մասնիկին՝ 1 (348), նախդիրներին՝ 6 (345-350), նա-
խադրասություններին՝ 6 (351-356) և կետադրական նշաններին՝ 4 (357-
360):

Հարկ է նշել, որ ընտրված քերականական կարգերի համախմբերի
թիվը գերազանցում է մինչև այժմ կատարված փորձերում ընդգրկված

Համախմբերի թվից: Մասնագիտական գրականության մեջ հայտնի է 120-150-ական կարգերի համախմբերի ընտրման դեպք: Մեր փորձի համար, ինչպես նշվեց, ընտրել ենք քերականական կարգերի թվով 360 համախմբը: Դա ունի և դրական, և բացասական կողմեր:

Դրականն այն է, որ տեքստի ամեն մի բառ նկարագրվում է ավելի մանրամասն, որը թույլ է տալիս հաշվի առնելու հնարավորին շափ ավելի շատ փաստեր: Այսպես, առանձին ուշադրության է արժանացել կետադրության հարցը, որը այլ փորձերում բաց է թողնվել (Բորոդկին): Մեր կարծիքով, տեքստում առկա կետադրությունը կարեոր ինֆորմացիա է կրում նրա կառուցվածքի վերաբերյալ: Տեքստի կետադրությունը կարելի է օգտագործել նախաղասությունների դասակարգման, նրանց սկզբնական և վերջնական հատվածների քննության, բառային կապակցությունների վերլուծության ժամանակ և այլն: Իհարկե, այն տեքստերում, որտեղ կետադրությունը բացակայում է, քերականական կարգերի ցանկից կարելի է հանել նաև կետադրությանը համարվող համախմբերը:

Քերականական կարգերի համախմբերի մեծ ցանկի բնտրության բացասական կողմն այն է, որ ստեղծվում են լրացուցիչ դժվարություններ զույգ հանդիպումների մատրիցի կազմման ու հետագա մշակման հարցում:

Որո՞նք են այդ դժվարությունների հաղթահարման միջոցները: Իհարկե, քերականական կարգերի 360 համախմբերի օգնությամբ տեքստի ձևաբանական վերլուծությունը և քերականական կարգերի զույգ հանդիպումների մատրիցի կառուցումը և մշակումը ավանդական եղանակներով շափաղանց դժվար է, գուցե և անհնար: Սակայն քերականական կարգերի այդպիսի քանակությունը ընտրվել է նկատի ունենալով աշխատանքների ավտոմատացումը EC-1045 հզոր հաշվողական մեքենայի օգնությամբ: Մյուս կողմից ծրագրային ապահովման համակարգը նախագծվել է այնպես, որ անհրաժեշտության դեպքում կարելի է նվազեցնել ելակետային կարգերի թիվը: Այսպես, եթե գոյականական կարգերի քննության ժամանակ հաշվի շառնվի հոլովումը, ապա 7 անդամ կպակասի գոյականին համադրվող քերականական կարգերի համախմբերի թիվը, և այն 168-ից կդառնա 24: Նման փորձեր կարելի է կատարել նաև մյուս խոսքի մասերի հետ: Պարզագույն դեպքում կարելի է փորձը դնել միայն խոսքամասային մակարդակով:

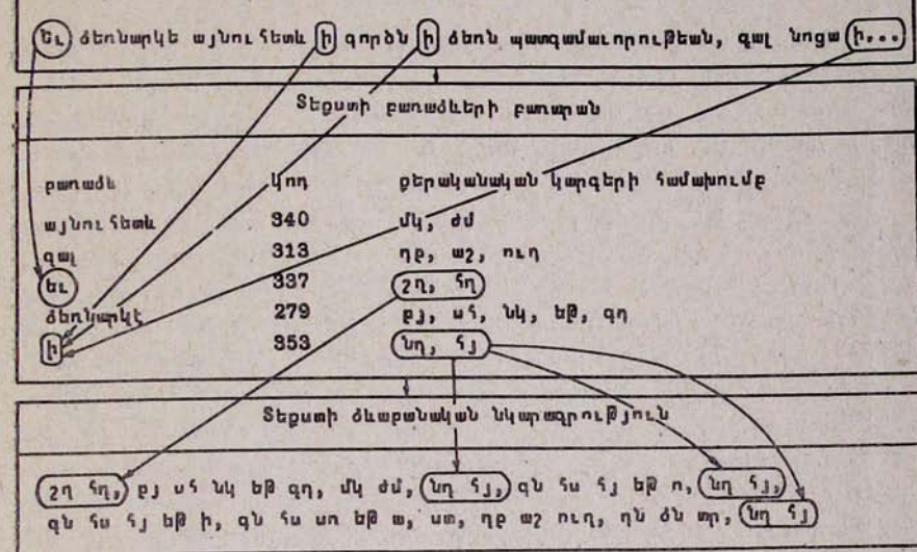
Այսպիսով, քերականական կարգերի զույգ հանդիպումների մատրիցի կազմման առաջին փուլը տեքստի ձևաբանական նկարագրությունն է, որը տեքստաբանության հիմնական, աշխատատար փուլերից է: Այդ նըկատառումով տեքստի վերլուծության ժամանակ առաջին հերթին դրվում է ձևաբանական վերլուծության ավտոմատացման խնդիրը: Ժամանակա-

Հից լեզուների համար արդեն գոյություն ունեն ձևաբանական վերլուծության ավտոմատ համակարգեր:

Հը. Աճապյանի անվան լեզվի ինստիտուտում մշակված է ժամանակակից հայերենի ձևաբանական վերլուծության ալգորիթմ (Ռ. Ռոռտյան): Ինչ վերաբերում է հին հայերենին, ապա համապատասխան ալգորիթմեր դեռ մշակված չեն: Այստեղ, եթե դիմենք ձևաբանական վերլուծության ավանդական նղանակներին, այսինքն՝ կատարենք ձեռքով, ապա խիստ նվազում է ընտրանքը, որը բացասաբար կանդրադառնարդյունների հավաստիության և ճշգրտության վրա:

Այդ պատճառով մենք ընտրել ենք տեքստի ձևաբանական նկարագրության ավտոմատացման եղանակ, որն այսօր էլ լայն կիրառում ունի տեքստերի ավտոմատ մշակման գործում (Լենդերս): Դա ձևաբանական վերլուծության բառարանային նղանակն է:

ՏԵՇՍ



Եթե կազմենք տեքստի բառաձևերի բառարան, որտեղ ամեն մի բառաձև համադրվում է բերականական կարգերի համախմբով, ապա այն կարող ենք օգտագործել ձևաբանական նպատակով: Դրա համար անհրաժեշտ է տեքստից ընտրած յուրաքանչյուր բառի համար բառարանում գտնել համապատասխան գլխաբառը և ստանալ նրան համադրվող քերականական կարգերի համախումբը:

Տեքստի ձևաբանական վերլուծության առաջարկվող եղանակը կարելի է ներկայացնել հետևյալ կերպ.

Որպես օրինակ բերենք Մ. Խորենացու «Հայոց պատմություն» ստեղծագործությունից (Մ. Խորենացի, Հայոց պատմություն, Երևան, 1981) 96էջի առաջին պարբերության ձևաբանական նկարագրությունը

Տեքստ Եւ ձեռնարկէ այնուհետև ի գործն ի ձեռն պատգամաւորութեան, գալ նոցա ի տեսութիւն միաբանութեան ի տեղի միջոց սահմանաց երկոցունց թագավորութեանցն. իբր բան ինչ և գործ հարկաւոր հասեալ, որ ի ձեռն գրոյ և հրեշտակութեան կատարել ոչ է կարողութիւն, եթե ոչ և դէմք երկոցունց հանդէպ լինիցին: Այլ գիտելով Տիգրանայ զառաքելոյ իրին կատարումն՝ ոչ ինչ յորոց խորհէքն Աժդահակ՝ ծածկէ, այլ ի ձեռն գրոյ յայտնէ որ ինչ ի նորայն խորութեան սրտի: Եւ յայտնեալ այսպիսոյ շարութեան՝ ոչ ինչ էր այնուհետև բան և խորամանկութիւն, որ զայսպիսի առագաստէր զշարութիւն. այլ յայտնի այնուհետև գրգոէր խաղմն:

Հատվածի ձեաբանական նկարագրություն. լշղ, հղլ բյ, սհ, նկ, եթ, դղ
/մկ, ժմ/ նդ, հյ/ զն, հս, հյ, եր, ո/ նդ, հյ լփն, հս, հյ, եթ, ի/ զն, հս, սո,
եթ, ա/ստ/ դբ, աշ, ուղ, լդն, ձն, տր/ նդ, հյ լփն, հս, հյ, եթ, ալ զն, հս, սո,
եթ, ա/նդ հյ/ զն, հս, հյ, եր, ո/զն, հս, սո, եթ, ո/զն, հս, սո, հյ, ալ
թն, ած, սո, հյ/ զն, հս, սո, հյ, ալ մտ նիս, հյ/ զն, հս, ուղ, եթ, ի լդն, աշ,
ուղ լշղ, հղլ զն, հս, ուղ, եթ, ո/ան, դբ, ուղլ դբ, աց, ուղ լստ/ զն, հբ, ուղ
նիս, սոլ զն, հս, սո, եթ, ո լշղ, հղլ զն, հս, սո, եթ, ա/դբ, աշ, ուղլ/ մս
լրյ, սհ, նկ, եթ, գդլ զն, հս, ուղ, եթ, ա/ստ/ լշղ, հդ լմսլ զն, հս, ուղ, հյ, ի
թն, ած, սո, եթ/ մկ, տղ լրյ, սդ, առ, հյ, գդլ վզ լշղ հդ լդբ, աշ, գծ/ զն,
հտ, սո, եթ, ա/նդ, հյ/ զն, հս, սո, եթ, ոլ զն, հս, սո, եթ, ա/զն, հս, հյ,
եթ, ալ րթ լմսլ զն, հբ, հց ննդ, բց/ զն, հբ, բց լրյ, սհ, աա, եթ, գդլ զն, հտ,
ուղ, եթ, ա/թթ/ բյ, սհ, նկ, եթ, գդ լստ/ լշղ, հդ ննդ, հյ/ զն, հս, հյ, եթ, իա
յզն, հս, սո, եթ, ոլ բյ, սհ, նկ, եթ, գդ լշղ, սդ լնդ, հց լփն, հբ, ուղ լնդ, տր/ զն,
ան, տր եթ, ալ զն, հս, տր, եթ, ի/ վզ լշղ, հդ լդբ, աց, ուղ լդն, ցց,
սոլ զն, հս, սո, եթ, ա/թթ/ մս լդն, հբ, ուղլ բյ, սհ, աա, եթ, գդ լմկ, ժմ
զն, հս, ուղ, եթ, ի լշղ, հղլ զն, հս, ուղ, եթ, ա/ստ/ զն, հբ, ուղ լզն, ցց, հյ/ բյ,
սհ, աա, եթ, գդ լնդ, հյ/ զն, հս, հյ, եթ, իա լմտ/ լշղ, հդ լան, դբ, ուղլ
մկ, ժմ լրյ, սհ, աա, եթ, գդլ զն, հս, ուղ, եթ, ի լվզ

Եթե տեհրստի ձևաբանական նկարագրությունը պատրաստ է, արդեն կարելի է անցնել քերականական կարգերի զույգ հանդիպումների մատրիցների կազմմանը: Մեր կողմից բերած պարբերության մեջ ներառված է 100 բառ, հաշված նաև կտտադրական նշանները:

Սովորաբար փորձարկվում են 1000-1500 լիիմաստ բառերով հատվածներ: 100 բառը նկարագրվում են քերականական կարգերի 46 համախմբով (քերականական կարգերի համախմբերի ցանկը տրված է հավելվածում): Հատվածի բառերի և քերականական համախմբերի հարաբերությունը 2:1 է: Սակայն բառերի հետագա ավելացումը չի բերի քերականական կարգերի համեմատական աճի, քանի որ համախմբերի թիվը 360 է: 1500 բառի դեպքում այդ հարաբերությունը կլինի մոտավորապես 5:1:

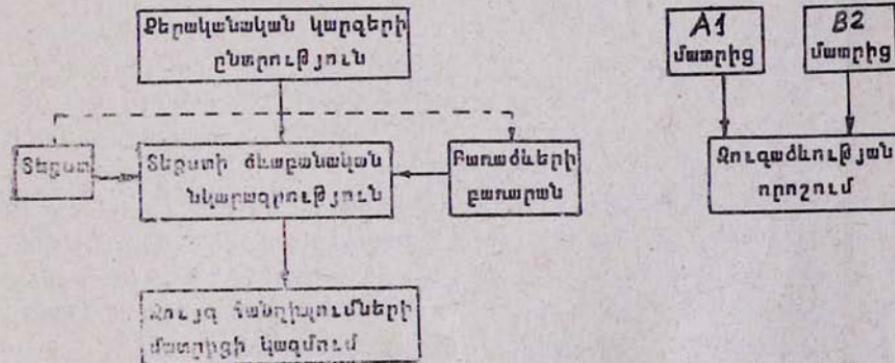
Բերված 100 բառն արդեն ցույց է տալիս քերականական կարգերի և նրանց դուգորդումների որոշ բաշխվածություն: Այն ավելի ակնառու դարձնելու համար զույգ հանդիպումների մատրիցի հիման վրա կառուցվում է զույգ հանդիպումների գրաֆ: Դա մի ուղղորդված հարթ գրաֆ է, որի գագաթները փոխհամապատասխանության մեջ են դրվում քերականական կարգերի համախմբերի, իսկ աղեղները, որոնք ուղղորդված են ըստ տեքստի ձախից աջ, ցույց են տալիս զույգ հանդիպումների թիվը:

Փորձնական եղանակով հաստատված է, որ նույն ոճի տեքստերին համապատասխանող գրաֆների զուգաձևության (հզօմօրֆոնտ) աստիճանը բարձր է: Զուգաձևության արդեն 25—30 տոկոսի դեպքում կարելի է խոռնել երկու տեքստերի միևնույն ոճի մասին:

Տեքստում քերականական կարգերի զույգ հանդիպման բաշխվածության խնդրի մաթեմատիկական ապահովման համակարգի սկզբունքային սխեման բերվում է ստորև (նկ. 1):

1

2



Համակարգը բաղկացած է երկու ենթահամակարգերից՝ 1) զույգ հանդիպումների մատրիցի կազմման և 2) զույգ հանդիպումների մատրիցների համեմատության:

Առաջին հնիտակամակարգի ելակետային տվյալները տեքստն է, աշխատանքի արդյունքը՝ քերականական կապերի զույգ հանդիպումների մատրիցը: Աշխատանքի սկզբունքը հետևյալն է. նախ բառաձենների բառարանի օգնությամբ ստացվում է տեքստի ձևարանական նկարագրությունը, որը, այնուհետև վերածվում է քերականական կարգերի համախմբերի կողերի հաջորդականության: Այսպես, եթե վերևում բերած պարբերության ձևարանական նկարագրության մեջ (շղ, հդ) քերականական կարգը փոխարինենք համապատասխան կողով՝ 337 (տես ծանոթությունը), բյ, սկ, նկ, եր, գդ, կարգը՝ 297-ով և այլն, ապա կստանանք կողերի հետևյալ հաջորդականությունը 337, 279, 340, 354, 130, 354, 129, 99, 359, 313, 231, 54 ...

Վերցնենք հաջորդականության առաջին զույգը՝ 337, 279:

Այն ցույց է տալիս, որ 337 և 279 կողերը ունեցող քերականական կարգերով բնորոշվող բառերը տեքստում հարադիր են, որտեղ 337 կողով բնորոշվող բառը զբաղեցնում է ձախադաս դիրք: Հստ զույգ հանդիպումների մատրիցի սահմանման 337, 269, կողերի զույգին համապատասխանության մեջ է դրվում զիյ=337, 279 տարրը, որովհետև, ինչպես ընդունել ենք, տեքստում մատրիցի տողերին համապատասխանող կարգերին հաջորդում են սյուներին համապատասխանող կարգերը: Հետեւարար, եթե վերևում ստացված հաջորդականության անդամները զույգ առ զույգ խմբավորենք այնպես, որ միևնույն կողը առաջին խմբում ունենա ձախադաս դիրք (ի-ի) տեղում, ապա կստանանք զույգ հանդիպումների մատրիցի տարրերը:

(337, 279), (279, 340), (340, 354), (354, 130), (130, 354), (354, 129), (129, 99), (99, 359), (395, 313), (313, 231), (231, 354)...

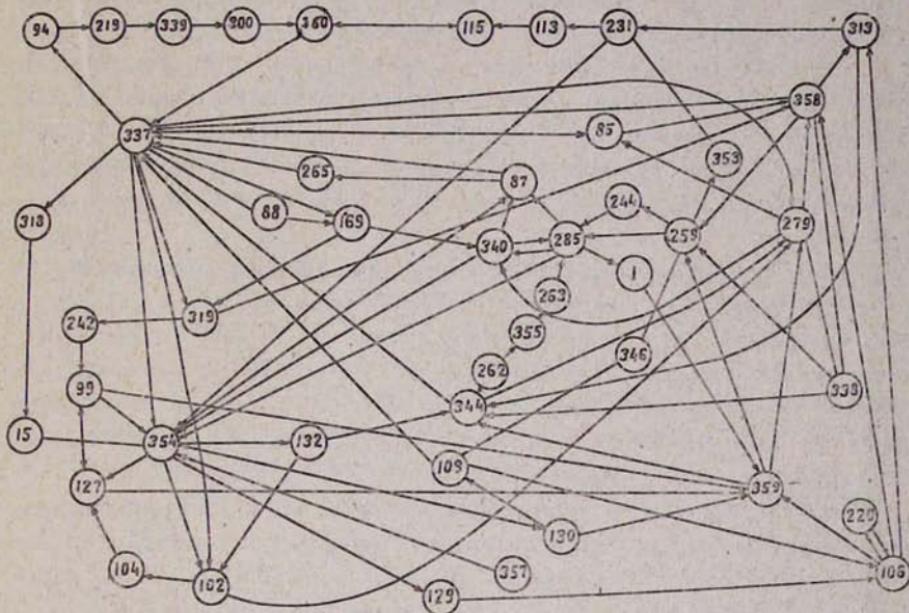
Ակնհայտ է, որ նման խմբավորման դեպքերում կլինեն խմբեր, որոնց բովանդակությունը կհամընկնեն: Դա ստացվում է այն դեպքում, երբ երկու հանդիպումներում համապատասխանարար ձախադաս և աջադաս բառերի քերականական կարգերը համընկնում են: Մեր կողմից քերկած օրինակներում (նդ, հյ) և (գն, հս, հյ, եթ, ի) կարգերը հանդիպում են երկու անգամ: Եվ քանի որ զույգ հանդիպումների մատրիցի տարրերը ցույց են տալիս տողերին ու սյուներին համապատասխանող կողերի (քերականական կարգերի) հանդիպումների թիվը, ապա հաջորդ փուլում ստացված հաջորդակցությունից պետք է հեռացնել կրկնվող խմբերը, նախապես գրանցելով նրանց հանդիպման թիվը.

(337, 279), (279, 340), (340, 354), (354, 130), (130, 354), (354, 129), (129, 99), (99, 359), (359, 313), (313, 231), (231, 354)...

Այժմ ստացված հաջորդականության հիման վրա զույգ հանդիպումների մատրիցի կազմումը առանձին դժվարություն չի ներկայացնում:

Զույգ հանդիպումների մատրիցի կաղման հարցում կարևոր խնդիր է տեքստի ձևալին նկարագրության ստացումը: Խնչակես նշվեց վերևում, այստեղ այդ խնդիրը վճռվում է տեքստի լեզվի բառաձևերի բառարանի օգնությամբ: Դա մի ցուցակ է, որի մուտքը բառաձևեր են, իսկ ելքը՝ դրանք բնորոշող քերականական կարգերի համախմբեր: Առանձին խնդիր է այդ ցուցակի կազմումը: Այն կարելի է վճռել երկու եղանակով.

ա) նախապես կազմվում է լեզվի ամբողջ բառապաշտի բառաձևերի բառարանը բ) կազմվում է հաշվողական մեքենայի մեջ գրանցված տեքստի բառաձևերի բառարանը: Մենք ընտրել ենք վերջինը:



Հարկ է նշել, որ եթե բառաձևերի բառարանը դիտենք իբրև առանձին հանգույց, որը կարելի օգտագործել տարբեր խնդիրների վճռման համար (տես ներքեւում), ապա, հատուկ մաթեմատիկական ապահովման դեպում, երկրորդ եղանակով կարող ենք հասնել ամբողջ բառապաշտի բառաձևերի բառարանի կազմմանը: Դրա համար անհրաժեշտ է նորանոր տեքստերի մշակման ժամանակ զուգահեռաբար լրացնել բառաձևերի բառարանը, մինչև ստացվի այնպիսի վիճակ, երբ այն այլևս լրացման կարիք չունենա:

Համակարգի մաթեմատիկական ապահովման առումով կարևոր հարց

է քերականական կարգերի ընդհանուր ցանկից տվյալ փորձի համար անհրաժեշտ համախմբերի ընտրությունը:

Քերականական կարգերի համախմբերի ընդհանուր թիվը 360 է: Սակայն քերականական կարգերի այդպիսի մեծ քանակ բոլորովին էլ պետք չէ ամեն մի փորձի համար, և փորձի կառավարումն ավելի ճկուն դարձնելու նպատակով մենք պետք է հնարավորություն ունենանք ընտրելու այն քերականական կարգերը, որոնք անհրաժեշտ են տվյալ տեքստի նկարագրության համար: Այլ կերպ ասած՝ համակարգը մեղ պետք է հնարավորություն տա թվարկել այն քերականական կարգերը, որոնց օգնությամբ կազմելու է տեքստի ձևաբանական նկարագրությունը: Այդ խնդիրը լուծվում է հատուկ երկխոսական ծրագրի օգնությամբ:

Եթե երկխոսության ընթացքում համակարգին որևէ լրացուցիչ ինֆորմացիա չի հաղորդվում, ապա տեքստի ձևաբանական վերլուծությունը կատարվում է այսպես կոչված խոսքիմասային մակարդակով՝ գոյական, ածական, թվական, դերանուն, բայց, դերբայց, շաղկապ, մակբայց, նախորդիր, նախադրություն, ժխտական մասնիկ, միջակետ, ստորակետ, բութ, վերջակետ:

Հետեւաբար, զույգ հանդիպումների մատրիցի նվազագույն չափը 16×16 է, այն դեպքում, երբ առավելագույնը՝ 360×360 է:

Այսպիսով, համակարգը մեղ թուզ է տալիս, ըստ անհրաժեշտության, տեքստի ձևաբանական վերլուծության խորության (մանրամասնության) աստիճանն ավելացնել ավելի քան 500 անգամ (եթե 16×16 մատրիցի դեպքում առավելագույնը կունենանք 256 հանդիպում ապա 360×360 մատրիցի դեպքում՝ 129600 հանդիպում):

Երկրորդ ենթահամակարգը նախատեսված է երկու մատրիցների դուգաձեռնության աստիճանի որոշման համար: Ստացված մատրիցները վերծելով երկուորդական մատրիցների՝ դիտելով յուրաքանչյուր հանդիպումը իրեն կապ տողերին և սյուներին համապատասխանող տարրերի միջև, զուգահեռաբար կարող ենք հաշվել այդ կապերի որոշ ընութագրիչներ՝ մեկ, երկու և ավելի կապերով հեռացված տարրերի մատրիցներ:

Անդրադասնանք երկրորդ խնդրին:

Ամեն մի տեքստի առանձնահատուկ է տեքստաշափական այնպիսի մի բնութագրիչ, ինչպիսին բառարանի ծավալն է:

Եթե համեմատենք լեզվի բառարանային միավորները տեքստային միավորների հետ, ապա նրանք տեքստում հանդես են գալիս անփոփոխ կամ փոփոխված՝ հարուցուցային (հոլովում, խոնարհում) շարքի որևէ անդամի տեսքով:

Հայտնի է, որ ամեն մի հարուցուց ոմնի իր ելակետային, կամ ինչպես ասվում է չնշույթավորված, միավորը, որին հակադրվում են շարքի մյուս՝

Նշուլթավորված միավորները: Տեքստի բառերի՝ որպես նշուլթավորված միավորների, ելակետային ձևերի ամբողջությունը կազմում է տեքստի բառարանը:

Վերջինի ընտրությունը որոշվում է հեղինակի անհատականությամբ. մի հեղինակ հակում ունի լավագույն կերպով կիրառել լեզվի բառափոխման (հոլովման, խոնարհման) համակարգը և փոքր քանակի բառերով կերառել տեքստ, իսկ մեկ այլ հեղինակ նույն երկարության տեքստ կերտում է համեմատարար մեծ բառարանի օգնությամբ:

Տեքստաշափության մեջ ընդունված է, այսպես կոչված, զանազանության գործակիցը, որը ցույց է տալիս, թե հեղինակն ինչ շափով է օգտվում լեզվի բառափոխման համակարգից: Դա տեքստի բառերի քանակի և նրա բառարանի ծավալի հարաբերությունն է:

Ձանազանության գործակիցը ստատիկ պարամետր է, որը միայն ցույց է տալիս տեքստը որքան անգամ է մեծ բառարանից:

Տեքստաշափության մեջ կիրառվում է մեկ այլ բնութագրի ևս, որը ցույց է տալիս նոր բառերի կախումը բառակիրառությունների թվից: Ակնհայտ է, որ այն գծային չէ:

Փորձնական եղանակով որոշված է, որ տեքստում նոր բառերի և բառակիրառությունների թիվը կապված են հետևյալ առնչությամբ՝

$$V = RN^a$$

որտեղ

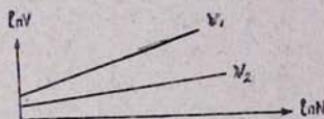
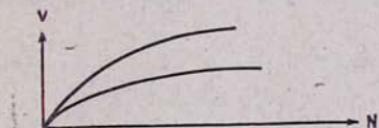
V -ն նոր բառերի թիվն է.

N -ը բառակիրառությունների թիվն է.

R -ը և a -ն տվյալ տեքստի համար բնորոշ հաստատուններ են ($0 < a < 1$):

Պարզ է, որ

R -ի և a -ի տարրեր արժեքների դեպքում կոտանանք տարրեր կորերի բազմություն:



Կորերի տեսքից պարզ երևում է, որ տեքստի սկզբնական մասում նոր բառերի թիվն արագ աճում է, իսկ այնուհետև հասնում մի սահմանի, որից աճը դադարում է:

Այդ գեղքում, երբ ընտրանքը մեծ է, հարմար է օգտվել լոգարիթմական կոռորդինատային համակարգից: Լոգարիթմելով վերևում տրված առնչությունը կստանանք.

$$I_n V = I_n R + a I_n N$$

Լոգարիթմական կոռորդինատային համակարգում սա ուղիղ գծի հավասարում է:

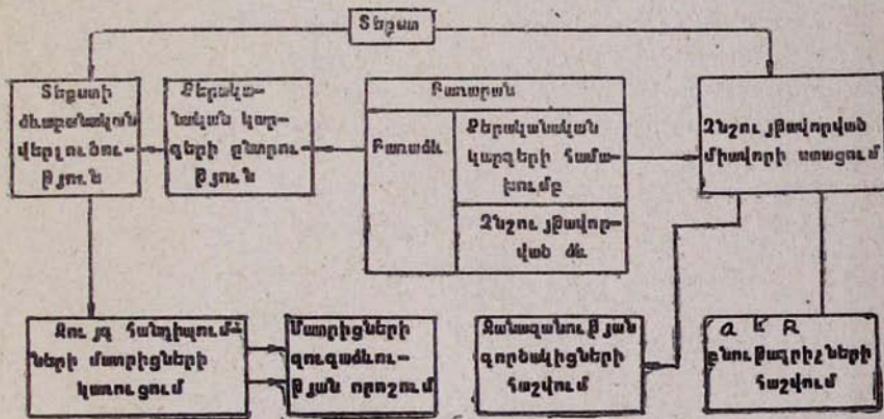
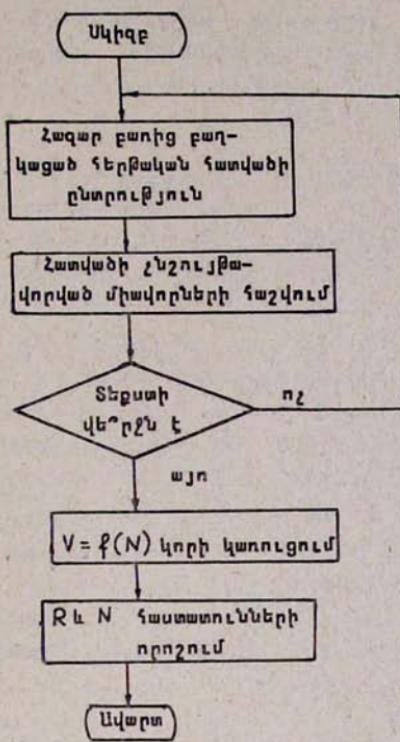
Այսպիսով, Մ. Խորենացու «Հայոց պատմություն» ստեղծագործության տեքստի համար զանազանության գործակցի և բառակիրառություններից կախված բառարանի ամի որոշման համար անհրաժեշտ է տեքստի ամեն մի բառը ձևափոխել չնշյունափոխված տեսքի: Դրա համար իրեն հիմք կարող է ծառայել բառաձեւերի բառարանը, որ կազմել ենք նախորդ խնդրի վճռման համար: Այդ բառարանի յուրաքանչյուր գլխաբառը լրացվում է նշույթավորված ձևով:

Ստորև բերում ենք բառաձեւերի բառարանից մի հատված, որտեղ գլխաբառերը լրացված են նշույթավորված ձևով:

Բառաձեւեր	Ըստարաւ
Նշույթավորված մեջ	Նշույթավորված մեջ
ոյնուհետեւ	ոյնուհետեւ
զալ	զալ
զրդն	զրծն
ձեռն	ձեռն
մեռնարկէ	մեռնարկէ
նոցա	նոցա
լատզամաւորութեան	լատզամաւորութեան

Ունենալով ալսպիսի բառաշան՝ ալգորիթմական եղանակով կարող ենք կազմել տեքստի չնշույթավորված ձևերի բառարանը: Որպեսզի ստանանք $V=RN^a$ բառաձեւերի R և գ հաստատունները, անհրաժեշտ է տեքստի ամեն մի հազարական հատվածի համար հաշվել նոր դոյացող չնշույթավորված ձևերի թիվը:

$V=RN^a$ կախման R և գ հաստատունների որոշման ալգորիթմը կարելի է ներկայացնել հետևյալ կերպ (տե՛ս ալգորիթմը՝ էջ 149): Յ-րդ նկարում տրվում է տեքստում քերականական կարգերի զույգ հանդիպումների մատրիցի կազմման ու վերլուծության, բառերի զանազանության գործակիցների հաշվման համակարգի սկզբունքային սխեման:



Սանոքություն. Տեքստում բերված հատվածի քերականական կարգերի զույգ հանդիպումների մատրիցի կազմման համար քերականական կարգերի ընդհանուր ցուցակից ընտրված համախմբերը համապատասխան կողմերով.

գն, հտ, ուղ, եթ, ա—1	դբ, աշ գծ—318
գն, հտ, սո, ա—15	դբ, աց, ուղ—319
գն, հս, ուղ, եթ, ա—85	շղ, հդ—337
գն, հս, ուղ, եթ, ի—87	շղ, սդ—338
գն, հս, ուղ, եթ, ո—88	մկ, տղ—339
գն, հս, ուղ, հթ, ի—94	մկ, ժմ—340
գն, հս, սո, սո, եթ, ա—99	մս—344
գն, հս, սո, եթ, ո—102	նխ, սո—346
գն, հս, սո, եթ, իս—104	նխ, հլ—348
գն, հս, սո, հթ, ա—106	գն, հս, հլ, եթ, իս—132
գն, հս, սո, հթ, ո—109	ան, գր, ուղ—169
գն, հս, տր, եթ, ա—113	թն, աձ, սո, եթ—219
գն, հս, տր, եթ, ի—115	թն, աձ, սո, հթ—220
գն, հս, հլ, եթ, ա—127	դն, ձն, տր—231
գն, հս, հլ, եթ, ի—129	դն, ցց, սո—242
գն, հս, հլ, եթ, ո—130	դն, ցց, հլ—244
դն, հր, հլ—262	դն, հր, ուղ—259
դն, հր, ցց—263	նդ, տր—353
դն, աշ, ուղ—265	նդ, հլ—354
բլ, սհ, նկ, եթ, գդ—279	նդ, ցց—355
բլ, սհ, աա, եթ, գդ—295	մտ—357
բլ, սդ, աո, հթ, գդ—300	ստ—359
դբ, աշ, ուղ—313	վշ—360

ՕԳՏԱԳՈՐԾՎԱԾ ԳՐԱԿԱՆՈՒԹՅԱՆ ՑԱՆԿ

- Бородкин Я. Н., Милов Л. В., О некоторых аспектах автоматизации текстологического исследования, в кн.: «Математические методы в историко-экономических и историко-культурных исследованиях». Наука, М., 1977, стр. 235—280.
- Виноградов В. В., Проблемы авторства и теория стилей, М., 1961.
- Мелихов А. Н. Ориентированные графы и конечные автоматы, М., Наука, 1971.
- Шмидт З. М. «Текст» и «История» как базовые категории в сб.: «Новое в зарубежной лингвистике», вып. VIII, М., 1978, стр. 89—110.
- Ուսուայան Ռ. Լ., Ժամանակակից հայերենի անվանական բառափոփոխման ձևալին նը-կարագրություն, «Հեղվի և ոճի հարցեր, Երդ պրակ» Երևան, 1982:
- Ուսուայան Ռ. Լ., Ժամանակակից հայերենի բայական ձևակազմություն, «Բառ, նա-խադասություն, տեքստ» Երևան, 1987:
- Maschinell Auswertung sprach-historischer Quellen. Ein Ber. Zur computertunterstützten Analys der Flexionsmorphologie des Frühneuhochchat (Lenciers W., Wegera K.—P., Berg E., et al.—Tübingen: Niemeyer, 1982—XI, 236 s.)