

LINGVISTIKA
LINGUISTICS

Tabiiy tilni qayta ishlash (nlp)da soʻz turkumlarini teglash masalasi

Manzura Abjalova¹
Dilrabo Elova²

Abstrakt

Kompyuter lingvistikasida tabiiy tilni qayta ishlash ancha murakkab jarayon boʻlib, unda ijtimoiy tabiatga ega tilning barcha hodisalari, umumiy va xususiy jihatlari, istisnoli holatlari, fonetik, morfologik, leksik, grammatik, semantik va hatto orfoepik xususiyatlarini eʼtiborga olish zarur hisoblanadi. Sunʼiy intellekt tizimini yaratishga bel bogʻlangan ayni damda matn birliklarini raqamli texnologiyalar orqali qayta ishlanishiga erishish muhim natija hisoblanadi. Bunday jarayonda soʻzlarning turkumini aniqlash zarur. Jahon kompyuter lingvistikasidan ushbu lingvo-texnik tahlilning soʻzlar turkumini aniqlash – PoS-tagger, yaʼni soʻz turkumlarini teglash, shuningdek, uning matnlarni avtomatik qayta ishlash jarayoni bosqichi ekanligi maʼlum. Til korpuslarini yaratish uchun boshlangan dastlabki teglash harakatlari bugungi kunga kelib, matn bilan bogʻliq koʻplab dolzarb masalalar yechimini bermoqda. NLP, yaʼni tabiiy tilni qayta ishlash jarayonida ham soʻz turkumlarini teglash birlamchi vazifa hisoblanib, buning natijasida omonimlikni, koʻp maʼnoli soʻzlar semantikasini aniqlash kabi turli lingvistik noaniqlik matn tarkibida tahlil qilinishiga erishiladi.

Mazkur maqolada soʻz turkumlarini teglash zarurati, matnlarni kompyuter tahlili jarayonidagi ahamiyati, teglash usullari haqida soʻz yuritiladi.

Kalit soʻzlar: *tabiiy tilni qayta ishlash, teglash, teg, soʻzlar turkumi, formal til, pragmatik xususiyat, korpus, polisemiya, omonimiya.*

¹Abjalova Manzura Abdurashetovna – filologiya fanlari boʻyicha falsafa doktori (PhD), Alisher Navoiy nomidagi Toshkent davlat oʻzbek tili va adabiyoti universiteti.

E-pochta: abjalovamanzura@navoiy-uni.uz

ORCID ID: 0000-0002-1927-2669

²Elova Dilrabo Qudratillayevna – oʻqituvchi, Alisher Navoiy nomidagi Toshkent davlat oʻzbek tili va adabiyoti universiteti.

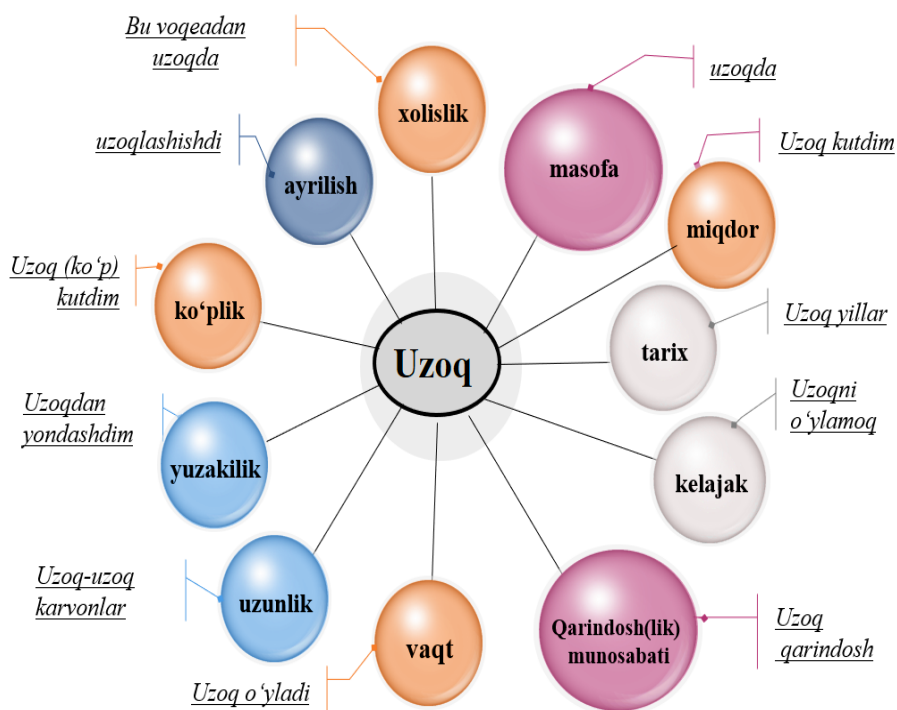
E-pochta: dilrabo@navoiy-uni.uz

ORCID ID: 0000-0002-1927-2968

Iqtibos uchun: Abjalova, M.A., Elova, D.Q. 2021. “Tabiiy tilni qayta ishlash (nlp)-da soʻz turkumlarini teglash masalasi”. *Ozbekiston: til va madaniyat* 1: 6–20.

Kirish

Inson borki, muloqotga oshiqadi, oʻz fikrlari va tuygʻularini izhor etishni xohlaydi. Axborot asrida boʻljak haqiqatni virtual tasavvur qiling: yoningizdagi robot itga “Men seni yaxshi koʻraman” (mehr, samimiy munosabat izhor etilyapti) va “Men seni yaxshi koʻryapman” (koʻzning yaxshi koʻrayotganligini anglatadi) jumllarini aytdingiz. Har ikki jumlani eshitib it dumini qimirlatadi. Lekin u bu jumllarni qanday maʼnoda qabul qilgan boʻladi? Robot it “Men seni yaxshi koʻraman, Qoravoy” deganingizni eshitganida, u “yaxshi koʻrmoq” ibora ekanligini anglashi mumkin. Koʻzlaringizdan, yuzingizda namoyon boʻlib turgan mehringiz orqali bu ifodaning tuygʻu ekanligini va unga qandaydir munosabat bildirishini tushunadi. Lekin unga qarata aynan koʻrish aʼzoyingiz (koʻz)ning shoh pardasi yanada ravshan koʻra boshlaganini bildirib aytilgan “Men seni yaxshi koʻryapman” (buyam sevinch va hayajon bilan) jumladagi birikmani toʻgʻri maʼnoda tushunmasligi mumkin. Mana shu misol robotni, yaʼni mashinani biz biladigan tilda muloqot qilishni oʻrgatish kerakligini koʻrsatadi. Yuqoridagi misoldan ayon boʻlganidek, muayyan soʻzning turli xil kontekstda ishlatilishidagi farqni tushunish juda muhim. Tilning pragmatik xususiyatiga oid ushbu fikrga “uzoq” soʻzi orqali yana bir misol keltiramiz (1-chizma):



1-chizma.

Mashinani tilga o'qitish uchun nima qilmoq kerak?

Ma'lumki, deyarli barcha mamlakat xalq ta'limi maskanlarida boshlang'ich sinflaridanoq o'quvchi gapni o'qiydi va undagi ot, sifat, son, fe'l, ravish, olmosh so'z turkumlarini aniqlaydi. Mana shu so'z turkumlari **tegi** hisoblanadi. So'z turkumlarini teglash (ingliz tilida bu *part-of-speech tagging* (*POS tagging* yoki *PoS tagging* yoxud *POST*), rus tilida *частеречная разметка* deyiladi) matnni avtomatik qayta ishlash bosqichi bo'lib, uning vazifasi matnda qo'llangan so'z (shakl)larning turkumi va grammatik xususiyatlarini aniqlash hisoblanadi. Shu vazifasi bilan POS-tagging matnni avtomatik tahlil qilishning dastlabki bosqichlaridan biri sanaladi.

Korpus lingvistikasida so'z turkumlarini teglash, grammatik kategoriyalarni teglash va so'zlarni toifalashda noaniqliklarni bartaraf etish uchun so'zni faqat uning lug'atdagi shakliga asoslanib emas, balki matn (jumla)dagi ifodasi bo'yicha uning turkumlik tegi va jumla (xatboshi, ibora)da boshqa so'zlar bilan birikish imkoniyatini hisobga olish muhim sanaladi. Gap bo'laklari teglarini identifikatsiyalash bir muncha qiyin jarayon. Sababi o'zbek tilidagi jamiki so'zlarni universal holda 12 turkum doirasida teglash imkoniyati yo'q. So'z uning jumla tarkibida reallashish holati va N-gramm [Abjalova 2020, 73-77] so'zlarning semantik valentligiga binoan polifunksional bo'lishi mumkin. Masalan: "Shifoxonaga bemorni keltirishdi" va "Shifoxonaga bemor odamni keltirishdi" jumlalardan birinchisida bemor so'zi turkumlik belgisi (kim? so'rog'iga javob berayotgan tushum kelishigidagi so'z)ga ko'ra ot turkumi, 2-jumlada esa (qanday? so'rog'iga javob beryapti) sifat turkumi vazifasidagi so'z hisoblanadi. O'zbek tili izohli lug'atida mavjud 11 000 o'zlashma so'zlardan 66 ta xuddi shunday polifunksional so'zlar aniqlandi [O'zbek tilining izohli lug'ati 2006].

Bir so'zning kontekstda ifodalanishiga ko'ra turli turkumga xoslanishi so'z turkumlari (ST) teglari uchun umumiy parametрни belgilash imkonini bermaydi. Bu holat korpus uchun ST teglarini qo'lda bajarish imkonsizligini ko'rsatadi. Shuningdek, yangi kontekstlarning yuzaga kelishi va tilda neologizmlarning paydo bo'lishi teglashtirish jarayonining davomiyligini ko'rsatadi. Shu bois til korpusida ST teglashtirishda mashinali kodlashtirishga tayaniladi. Rus tilidagi matnlarni qayta ishlashda so'z turkumlarini teglashtirish uchun *Yandex* qidiruv tizimida *Mystem* morfologik analizatori, *Tree Tagger*'da rus tili *utuliti* [Abjalova 2020,157], Python dasturlash tili-da yaratilgan NLTK dasturiy kutubxona mavjud.

Soʻz turkumlarini teglash tarixi

PoS-tagging boʻyicha tadqiqotlar korpus tilshunosligi bilan chambarchas bogʻliq. Matnni kompyuterda tahlil qilish uchun ingliz tilining birinchi yirik korpusi 1960-yillarning oʻrtalarida Genri Kucher va V. Nelson Frensis tomonidan Braun Universitetida ishlab chiqilgan. Shu bois mazkur til korpusi “Braun korpusi” nomiga ega boʻlgan.

Brown korpusida soʻz turkumlarini teglash uchun koʻp yillar davomida inson tomonidan soʻzlar va ularning turkumlari roʻyxat qilingan. Qoidalarga asoslangan dastlabki harakatlar Green va Rubin dasturi yordamida qoʻlda tuzilgan yirik roʻyxat asosida amalga oshirildi. Bunda, asosan, lingvistik birliklarning ketma-ket doimiy joylashuvi eʼtiborga olingan. Masalan, artikldan keyin ot turkumi (noun) keladi, feʼl (verb) turkumi emas. Bunday qatʼiy qoidalarga asoslangan dastur 70% toʻgʻri ishladi. Uning natijalari bir necha bor qoʻlda tekshirildi va tuzatildi, keyinchalik foydalanuvchilar tomonidan ham tuzatishlarni yuborish imkoniyati yaratildi, natijada 70-yillarning oxiriga kelib teglashtirish (yoxud markirovka) deyarli mukammal boʻldi.

Ushbu korpusdan son-sanoqsiz soʻzlar chastotasi va ularning turkumlarini oʻrganish uchun foydalanilgan va boshqa koʻplab tillarda shunga oʻxshash teglarni shakllantirish hamda rivojlantirishga sababchi boʻlgan. Yaxshi teglangan taʼlimiy korpuslar til modelini sinovdan oʻtkazish va takomillashtirish uchun qimmatli manbadir. Matnlar korpusi tilshunosga grammatik qoidalarni ishlab chiqishda eʼtibordan chetda qolgan lingvistik va nutqiy vaziyatlarni koʻrsatib (yoxud eslatib) turadigan tabiiy til manbasidir.

Bir muncha vaqt soʻz turkumlarini teglash tabiiy tilni qayta ishlash (NLP)ning ajralmas qismi deb hisoblangan, sababi ayrim holatlar mavjudki, soʻzning turkumini uning (ayniqsa, polifunksional va koʻp maʼnoli soʻzlarning) semantikasini, hattoki kontekst pragmatikasini tahlil qilmasdan aniqlab boʻlmaydi. Kompyuter dasturlarida esa ularning mukammalligini taʼminlovchi semantik va pragmatik tahlil bosqichlarini yaratish ancha murakkab jarayon hisoblanadi.

Soʻz turkumlarini teglash zarurati

ST teglashtirish usullarini koʻrib chiqishdan oldin, avvalo, STni teglash nima uchun kerak va ulardan qayerda foydalanilishiga toʻxtalamiz (2-chizma).



2-chizma.

Eng muhimi, ST teglari tabiiy tilni qayta ishlash (Natural Language Processing / NLP) uchun eng birlamchi zaruriy lingvistik element hisoblanadi, shu bois STni teglash NLPda turli xil muammolarni soddalashtirish uchun dastlabki shart sifatida amalga oshiriladi.

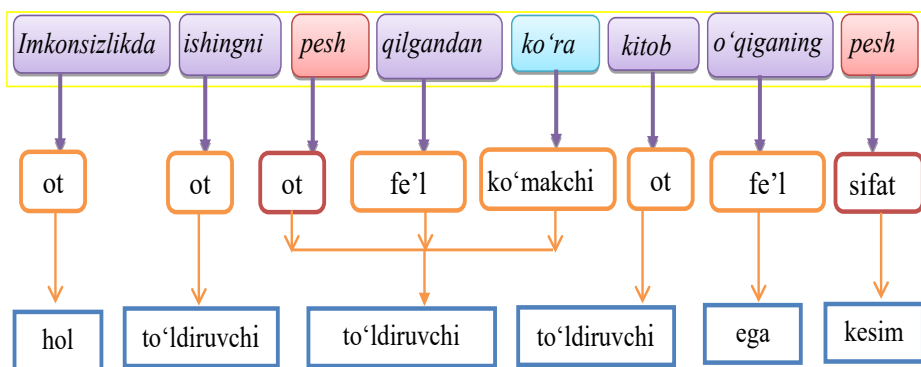
Tilning fonetik, leksik, morfologik, sintaktik va semantik sathlarini oʻzida qamragan universal lingvistik axborot tizimi (korpus)da aynan morfotahlilni amalga oshirish uchun ham soʻz turkumlarini teglash zarur hisoblanadi. Shunda istalgan soʻz (shakl)-ga oʻquvchi uchun toʻgʻri tahlil natijasi til korpusida taqdim etish imkoniyati boʻladi. Til korpusi mavjud boʻlmagan holda koʻpgina tillarda soʻzlarning turlanish va tuslanish jarayonlarida lingvistik taʼminotga kiritilgan morfologik lugʻatlarga murojaat qilinadi [Abjalova 2020, 38-39; 155].

Oʻzbek tilida 50 000 dan ortiq leksemalar mavjud boʻlib, korpus va lingvistik kompyuter dasturlari bazasida har bir leksemaning turkumi aniqlab berilishi muhim masala hisoblanadi. Shunday soʻzlar mavjudki, turkumlik belgisini oʻzida namoyon qilmaydi yoxud gap tarkibida kontekstual maʼnosi oʻquvchini chalgitadi. Masalan, "... test sinovlaridan oʻtkazish yuzasidan shaxsan javobgarligi belgilab qoʻyilsin" (<https://lex.uz/acts/-4276890>), "Shaxsan oʻzim keldim", "Shaxsan bajardim", "Bular hammasi lotincha yoki lotincha-ga yaqin soʻzlar. Men, shaxsan, shunday deb bilaman" (A.Qahhor. Adabiyot muallimi) jumllaridagi *shaxsan* yasama soʻzining turkumini aniqlash mushkul. Baʼzi oʻrinlarda oʻzlik olmoshi oʻrnida qoʻllanilayotgan leksema (olmosh), baʼzi hollarda ravish leksema tarzida namoyon etadi. Bu holatda soʻzning turkumini aniqlashda turkumlarning kategorial xususiyatlariga murojaat etiladi. Ular 4 ta: semantik, sintaktik, morfologik va soʻz yasalishi xususiyatlaridir.

Maʼlumki, oʻzbek tilida 12 soʻz turkumi (mustaqil soʻz turkumlari: ot, feʼl, sifat, ravish, son, olmosh; yordamchi soʻz turkumlari: bogʻlovchi, koʻmakchi, yuklama; alohida olingan soʻzlar turkumi:

modal, taqlid, undov)ga soʻz yasovchi qoʻshimchalarning qoʻshilishi natijasida 4 soʻz turkumi: ot, feʼl, sifat, ravish yasaladi. Aniqlangan yasovchi qoʻshimchalar (337 ta: ot yasovchi qoʻshimcha 114 ta, feʼl yasovchi 58 ta, sifat yasovchi 117 ta, ravish yasovchi qoʻshimcha 48 ta) [Abjalova 2020,122-123] sirasida **-an** ravish yasovchi affiks hisoblanadi. Ushbu parametrdan kelib chiqib xulosalash mumkinki, ot turkumiga mansub “shaxs” soʻziga **-an** yasovchi qoʻshimchasi birikishi natijasida yasama ravish hosil qilingan: shaxs (Ot) ~ {-an} => shaxsan.

Umuman olganda, soʻz oʻta murakkab hodisa sifatida talqin qilinib, ayni vaqtda ham til birligi, ham nutq birligi boʻlishi taʼkidlanadi. Til birligi bilan nutq birligining teng kelib qolishi holati, asosan, oʻzgaraydigan turkumlarga kuzatiladi. Bizga maʼlum boʻlgan poli-funksional, koʻp maʼnoli va omonim soʻzlar tabiati esa jiddiy tadqiqot va amaliy kuzatuv jarayonlarini talab etadi.



3-chizma.

“Imkonsizlikda ishingni pesh qilgandan koʻra kitob oʻqiganing pesh” gapida (3-chizma) soʻzlarning turkumlarga teglanishi natijasida gap boʻlaklarining aniqlanishi va “pesh” omonim soʻzining ketma-ket kelgan soʻzlar (pesh qilmoq (Ot + F); oʻqiganing pesh (F + Sif) va joylashuv oʻrniga koʻra (gap soʻnggida ot-kesim => Sifat) turkumi aniqlanishini koʻrish mumkin.

Soʻz turkumi har qanday soʻz birikmasi (3-chizmada “pesh qilgandan koʻra” – vositali toʻldiruvchi uch turkumdagi soʻzlar birikuvidan iborat) va gap strukturasi turli sintaktik funksiyalarni bajarishi mumkin. Har bir soʻz turkumida asosiy, birlamchi sintaktik funksiya mavjud boʻladi. Birlamchi sintaktik funksiya soʻz turkumining leksik maʼnosidan kelib chiqadi va bu maʼnoning oʻziga xos transpozitsiyasi sifatida gavdalanadi [Ganiyeva 2019, 277].

ST teglash uchun lingvistik bazada soʻzlar va ularning turkumlari koʻrsatilgan roʻyxatning kiritilishi kifoya emas. Yuqoridagi soʻz turkumini aniqlash holatidagi kabi izchillikning yoʻqolishi yoxud bir shaklga ega polifunksional, omonim yoki koʻp manoli soʻzlarning gapda ifodalagan turkumini topish hatto mutaxassis tilshunosni ham fikr yuritishga, izlanishga undaydi. Shuningdek, oʻzbek tilidagi koʻpgina soʻzlar muayyan turkumga mansubligi aniqlanmagan. Har bir tabiiy tilda mavjud bunday muammolar eʼtiborga olinib STni teglashda bir necha usullarga tayaniladi.

Soʻz turkumlarini teglash usullari

Aksariyat hollarda soʻz turkumlarini teglashda quyidagi usul (metod, algoritm)larga asoslaniladi:

- qoidalarga asoslangan usul.
- stoxastik (yoxud statistik) usul.

Qoidalarga asoslangan PoS teglar.

Eng azaliy teglash usullaridan biri bu qoidalarga asoslangan POS-teglash sanaladi. Bunda, asosan, Brill usuli qoʻl keladi [Brill 1992]. Qoidalarga asoslangan teggerlar har bir soʻzni teglash uchun lugʻat yoxud leksikadan foydalanadilar. Agar soʻzda (polifunksional, omonim, koʻp maʼnoli soʻzlar nazarda tutilmoqda) bir nechta teglar boʻlsa, unda qoidalarga asoslangan teggerlar gapdagi soʻzning turkumlik tegini toʻgʻri aniqlash uchun qoʻlda yozilgan qoidalaridan foydalanadi. Yanada aniq teglarni berishda soʻzning lingvistik xususiyatlarini undan oldingi va keyingi soʻzlarni tahlil qilish orqali qoidalarga asoslanib belgilash orqali ham amalga oshirilishi mumkin. Masalan, qaratqich kelishigidagi ismga mansub soʻzdan soʻng kelgan lingvistik birlik egalik qoʻshimchasini olgan ot turkumidagi soʻz hisoblanadi. Masalan, mening kitobim, akamning uyi, Salimaning koʻylagi kabi. Demak, bunday holda soʻzning ot turkumligi oʻzining oldida kelayotgan qaratqich kelishigidagi ism orqali belgilanadi. Ingliz tilidan misolga eʼtibor qaratamiz: oldingi soʻz artikl boʻlsa, u holda undan keyin kelayotgan soʻz ot turkumiga oid leksik birlik sanaladi. Masalan, *an egg, a book, the train, the windows* kabi.

PoS teglaridagi bunday holatlar qoidalar shaklida kodlanadi.

Ushbu qoidalar quyidagilarni tashkil etishi mumkin:

1. Lingvistik meʼyorlarga asoslangan qoidalar. Tilning orfografik qoidalariga asoslangan yuzlab qoidalar umumiy, xususiy va istisnoli qoidalar bazasi tarzida shakllantiriladi [Abjalova 2020].

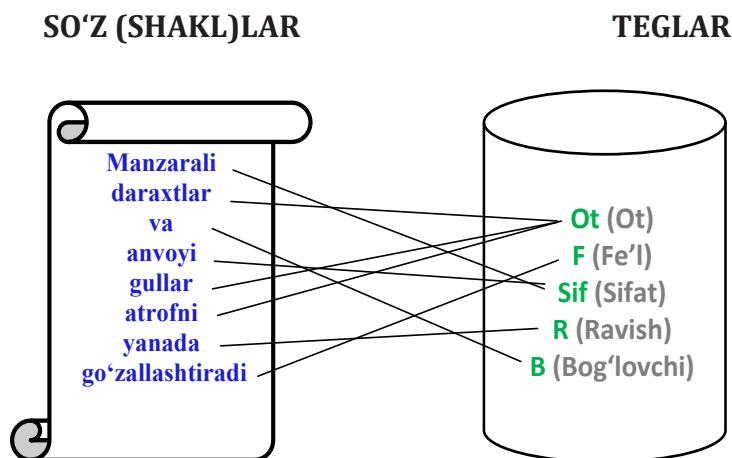
2. Kontekstli shablon qoidalari, yaʼni gap tarkibidagi koʻchma maʼnoga ega soʻzning muntazam ravishda konnotativ maʼnoda

qoʻllanishi dastur xotirasida saqlanadi, natijada keyingi jarayonlarda oʻsha konnotativ soʻz bilan bogʻliq noaniqliklar bartaraf etiladi.

Qoidalarga asoslangan metodga binoan soʻz turkumlarini teglash ikki bosqichda amalga oshiriladi:

Birinchi bosqichda PoS-tegger (izohli, imlo, morfologik yoxud orfografik) lugʻat(lar)ga tayanadi. U lugʻat yordamida har bir soʻzning turkum(lar)ini aniqlab chiqadi.

Ikkinchi bosqichda polifunksional yoki omonim soʻzlarning turkumlari qoʻlda yoziladi va bunday soʻzlarni jumla tarkibidagi vazifasini aniqlash qoidalarining katta roʻyxati ishlab chiqiladi.



4-chizma. Gapdagi har bir soʻzga soʻz turkumining belgilanishi.

Qoidalarini avtomatik generatsiya qilish metodining yaqqol namunasi – amerikalik lingvist Erik Brill metodidir [Brill 1995]. Ish uslubi quyidagicha:

1. Boshlash: Har bir soʻz oʻsha soʻzning tez-tez qoʻllaniladigan tegi bilan bogʻlanishi lozim. Nomaʼlum soʻzlar ot soʻz turkumi sifatida qabul qilinadi. Mazkur bosqichdan nafaqat oʻrganish jarayoni boshlanadi, balki omonimlarni bartaraf etish metodi ham harakatga keladi.

2. Tez-tez uchraydigan xatolik uchun oʻzgarish (qayta ishlash) qoidasini yaratish kerak.

3. Istalgan minimal xatoga erishguncha ikkinchi bosqichni takrorlash lozim.

Oʻtkazish (qayta ishlash) qoidalarida “eski teg, yangi teg, shart” toʻplamlari oʻz ifodasini topadi va qoidada belgilangan shart bajarilganda eski teg yangisi bilan almashtiriladi. Ushbu metodning kamchiligi qoidalar sonining koʻpayishi bilan aniqlik daraja-

sining pasayishida kuzatiladi [Brill 1992], bu Paretoning prinsipiga mos keladi: 80% harakat 20% natijani ta'minlaydi. Shu bilan birga, prinsip aksincha jarayonda ham ishlaydi: boshlash bosqichining faqat bitta qadamini bajarish natijasida omonimlikni bartaraf etishda yuqori aniqlikka erishish mumkin. SinTagRus korpusida o'tkazilgan test natijalaridan ma'lum bo'ldiki, ushbu metod har bir so'zning gapdagi o'rnini 97,4%, morfologik xususiyatlarning to'liq majmuini 87,6% aniqlik bilan topish imkonini beradi [Abjalova 2020,73].

Qoidalar asoslangan ST teglari xususiyatlari

Qoidalar asoslangan PoS teglari quyidagi xususiyatlarga ega:

- Ushbu teglar bilimga asoslanadi.

- Qoidalar qo'lda yaratiladi.

- Axborotlar qoidalar shaklida kodlanadi.

- Qoidalar cheklangan bo'ladi. Raqamli texnologiya uchun cheksizlik mavhumlikni ifodalaydi, "va shu kabilar", "... (ko'p nuqta), "va hokazo" kabi birikma va belgilar ro'yxatning davomiyligini emas, balki noaniqligini bildiradi. Shu bois kompyuter lingvistikasida bu kabi noaniqliklarga yo'l qo'yilmaydi, balki har bir lingvistik birlikning tabiati, xususiyatlari dasturiy ta'minot bazasida aniq ko'rsatilishi talab etiladi.

- Tilni modellashtirishda teglashtirish qoidasiga asoslaniladi.

- Stoxastik teglash usuli.

Mazkur usul chastota yoki ehtimollik (statistika)ka asoslanadi. Shu bois ayrim manbalarda statistik yoxud ehtimolikka asoslangan usul tarzida tushuntiriladi. Ayon bo'lganidek, oddiy stoxastik teglashda ST teglari uchun quyidagi metodlardan foydalaniladi:

1. Chastotali yondashuv.

Ushbu yondashuvda stoxastik teglar so'zning matnda ma'lum bir teg bilan uchrashi ehtimoli asosida grammatik noaniqliklarni bartaraf etadi. Shuni ham aytish mumkinki, o'rganilayotgan to'plam (matn qismi)da muayyan so'z bilan tez-tez qo'llaniladigan teg o'sha so'zning noaniqligiga ma'lumot berishga yordamchi tegdir. Masalan, qo'llanilish darajasiga binoan tushum kelishigidagi so'zdan so'ng kelgan lingvistik birlik fe'l turkumiga mansub hisoblanadi: kitobni o'qimoq, ismni yozmoq, mamlakatni aylanib chiqmoq => {Wni ~ V} kabi. Teglash jarayonidagi bunday yondashuvning asosiy muammosi, u tabiiy tilda birikuvchanligi bo'lmagan teglar ketma-ketligini keltirib chiqarishi mumkin. Masalan, "Do'stinga kitobni sovg'a qilish uchun sotib oldim" misolida sovg'a so'zi o'ng tomondan

“qilish uchun” birikmasi bilan birikuvchanlikka ega, ammo chapdagi soʻz bilan (kitobni) na grammatik, na semantik jihatdan birika oladi. Bunday hollarda teglashtirishda birikuvchanlikka ega boʻlmagan soʻzlar nomutanosibligiga asoslanish natijasida morfoanaliz jarayonida notoʻgʻri maʼlumot yuzaga keladi.

2. Teglarining ketma-ketligi ehtimoli yoxud n-gramma usuli.

Stoxastik usulning mazkur yondashuvi tegger berilgan teglar ketma-ketligining qoʻllanilish ehtimolini hisoblaydi. Ketma-ketlik oʻlchovi, yaʼni n (bigram – ikki element ketma-ketligi, trigram – uch ketma-ket teg, 4 gram – toʻrt teg ketma-ketligi) teglarga asoslangani uchun bu yondashuv N-gramma usuli ham deyiladi. N-gramma – matnlarga avtomatik ishlov berishda keng qoʻllaniladigan matematik hisob vositasidir. Oʻzbek kompyuter lingvistikasida S.Rizayev harf birikmalarini *bigramm*, *trigramm* terminlari bilan ifodalagan [Rizayev 2006,18].

Yashirin Markov modeli stoxastik usulda faol qoʻllaniladi. 1960-yillarda L.E. Baum va uning hamkasblari tomonidan ishlab chiqilgan [Baum, Sell 1968, 211-227] mazkur usul statistik jarayonda yuzaga keladigan barcha variantlar ehtimolligini hisobga olishga yordam beradi. Masalan, maʼlum bir matnda ot turkumiga oid soʻzlar bogʻlovchiga nisbatan tez-tez va koʻp uchrasa unda ayni kontekstda mavjud omonim katta ehtimollik bilan bogʻlovchi emas, ot turkumiga oid soʻz boʻladi, keyingi ehtimollikda bogʻlovchi sifatida hisobga olinadi. Kontekstni tavsiflash uchun N-grammadan foydalaniladi. N-gramma – soʻzlar yoki teglar kabi N-identifikator elementlarning ketma-ketligini ifodalaydi.

Yashirin Markov modellari termodinamika, statistik mexanika, fizika, kimyo, iqtisodiyot, moliya, signallarni qayta ishlash, axborot nazariyasi, nutqni qayta ishlash, husnixat, imo-ishoralarni tanib olish, [Starner, Pentland 1995] soʻz turkumlarini teglash va bioinformatikada keng qoʻllaniladigan statistik model hisoblanadi [Li 2003; Ernst 2012].

Soʻz turkumlarini stoxastik teglash usuli xususiyatlari

Stoxastik PoS-tegerlar quyidagi xususiyatlarga ega:

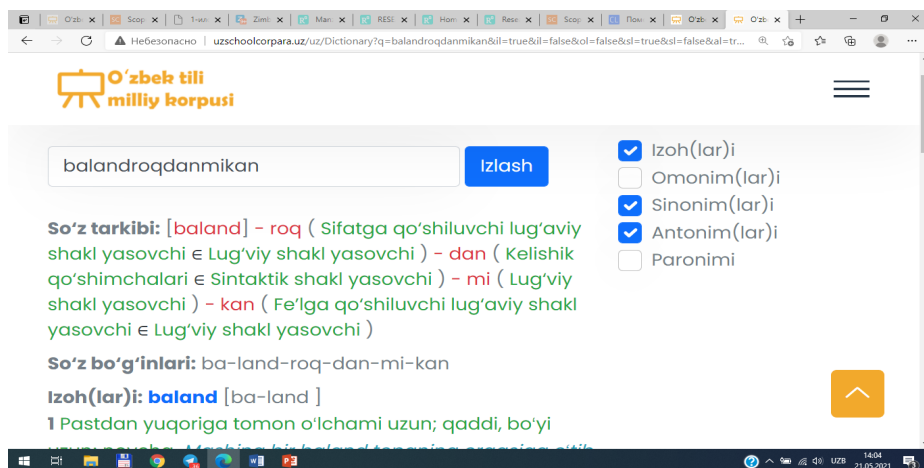
- Mazkur teglashtirish teglarning ketma-ket qoʻllanish darajasi boʻyicha ehtimolligiga asoslanadi.
- Taʼlim korpusi talab qilinadi.
- Korpusda mavjud boʻlmagan soʻzlar uchun hech qanday ehtimollik boʻlmaydi.
- Taʼlim korpusidan tashqari boshqa tur til korpusidan ham

foydalanish mumkin.

– Eng oddiy ST teglash usuli, chunki bu usulda til korpusidagi faol tarzda ketma-ket qoʻllangan teglarni tanlab oladi.

Amaliy natija

Bir necha yillar davomida olib borilgan tadqiqotlar va 2020-2021-yillardagi amaliy sa'y-harakatlar natijasida Toshkent davlat o'zbek tili va adabiyoti universiteti Axborot texnologiyalari hamda Amaliy tilshunoslik va lingvodidaktika kafedralari hamkorligida AM-FZ-201908172 raqamli "O'zbek tili ta'limiy korpusini yaratish" mavzusidagi amaliy loyiha doirasida O'ZBEK TILI MILLIY va TA'LIMIY KORPUSLARI yaratildi. Bugungi kunda mazkur korpusda mosfoanalizator (avtomatik morfologik tahlil), sinonimizator (qidiruvga yozilgan so'zga uning ma'nodoshlarini taqdim etish dasturi), shuningdek, so'zni bo'g'inlarga ajratish, izoh (lar)ini taqdim etish va antonimlarini ko'rsatish imkoniyatlari yaratilgan.

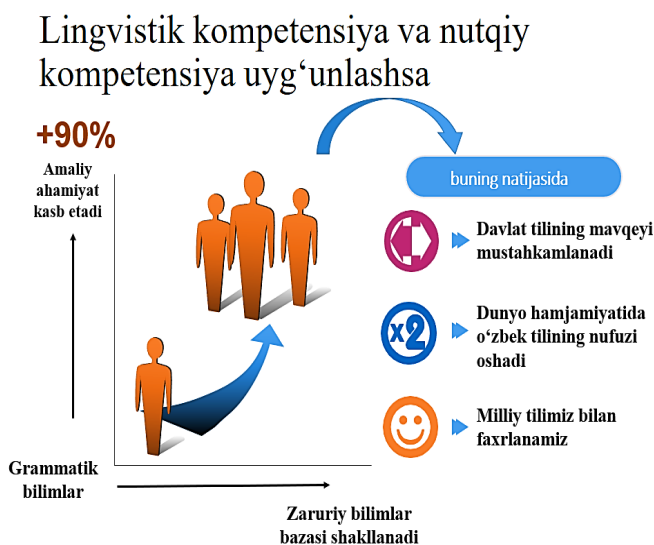


5-rasm. O'zbek tili Milliy korpusida "balandroqdanmikan" so'z shaklining morfologik tahlil ("So'z tarkibi" qismida) natijasi.

Xulosa tarzida aytish joizki, tabiiy tilni qayta ishlash (NLP)-da so'z turkumlarini teglash algoritmlarni yaratish avtomatik tahlil, morfoanaliz va tarjimon dasturlari uchun birlamchi shart amalga oshirilgan sanaladi, natijada matnlarning grammatik jihatdan sifatli tahlil qilinishiga erishiladi. Bugungi kunda neyrotarmoqqa asoslangan sun'iy intellekt tizimida qoidalarga asoslangan va stoxastik usullardan gibrid tarzda foydalanilmoqda.

Lingvistik bilimlar kompyuter dasturlari va til korpuslari bazasi uchun eng zarur ma'lumot manbai hisoblanadi. Lingvistik

prosessorda lisoniy maʼlumotlar bazasini shakllantirishda, kundalik hayotda til bilimlaridan foydalanish koʻnikmalariga asoslanilsa, bunday vaziyatlar adabiy til meʼyorlariga aylantirilsa, kompyuter dasturlarining tahlil imkoniyati mutaxassis darajasida mukammallikka yetadi. Shu bois lingvistik kompetensiya va nutqiy kompetensiya doimiy ravishda bir-birini taqozo etadi. Aynan soʻz (asosan, polifunksional, koʻp maʼnoli, omonim soʻzlar)ning turkumini belgilashda ham maqola avvalida “uzoq” soʻzi misolidagi kabi soʻzning pragmatik va kontekstual maʼnolari hamda qoʻshimcha semalariga asoslanilsa, raqamli texnologiya dasturlari va tizimlarining amaliy ahamiyati yanada oshadi (6-rasm).



6-rasm.

Umuman olganda, bugungi kunda lingvistik bilimlar tildan amalda foydalanish imkoniyatiga asoslanib mukammallashtirilsa, nazariy manbalarga tayanilib, kompyuter dasturlari va tizimlari uchun formal oʻzbek tili yaratilsa, kelajakda barcha turdagi (oʻzbek tarjimon dasturi, nutqni tanish, oʻzbek tili morfoanalizatori, kompyuter muloqoti kabi) dastur hamda elektron tizimlarining yuzaga kelishiga zamin yaratiladi.

Adabiyotlar

- Abjalova, M. 2020. *Tahrir va tahlil dasturlarining lingvistik modullari*. Toshkent: Nodirabegim.
- Abjalova, M., Yuldashev, A. 2021. "Methods for Determining Homonyms in Linguistic Systems". *ACADEMICIA: An International Multidisciplinary*

Research Journal 11 (2): 700-715. DOI: 10.5958/2249-7137.2021.00522.X.

- Abjalova, M. 2021. "O'zbek tili Milliy korpusida so'zshakllarni leksikografik baza asosida qidiruv imkoniyatlari". *Kompyuter lingvistikasi: muammo, yechim, istiqbollar*, 12-17. Toshkent: ToshDO'TAU.
- Abjalova, M. 2021. "O'zbek tili milliy korpusida sinonimayzer yoxud sinonimizatorni yaratish masalasi". *O'zbek Milliy va ta'limiy korpuslarining yaratishning nazariy hamda amaliy masalalari*, 38-40. Toshkent: ToshDO'TAU.
- O'zbek tilining izohli lug'ati*. 2006. 80 000 dan ortiq so'z va so'z birikmasi. A. Madvaliyev tahriri ostida. 5 jildli. Toshkent.
- Brill, E. 1995. "Transformation-Based Error-Driven Learning and Natural Language Processing: A Case Study in Part-of-Speech Tagging". *Computational Linguistics* 21: 543-565. <http://acl.ldc.upenn.edu/J/J95/J95-4004.pdf>.
- Brill, E. 1992. "A simple rule-based part of speech tagger". *Proceedings of ANLC* 154.
- Baum, L. E.; Sell, G. R. 1968. "Growth transformations for functions on manifolds". *Pacific Journal of Mathematics* 27 (2): 211-227.
- Ganiyeva, Dildora. 2019. "Мазмуний синкретизм ва полифункционаллик". *NamDU ilmiy axborotnomasi* 6: 275-278.
- Rizayev, S. 2006. *O'zbek tilshunosligida lingvostatistika asoslari*. Toshkent: Fan.
- Starner, Thad, Pentland, Alex. 1995. *Real-Time American Sign Language Visual Recognition From Video Using Hidden Markov Models*. Master's Thesis, MIT, Program in Media Arts.
- Li, N; Stephens, M. 2003. "Modeling linkage disequilibrium and identifying recombination hotspots using single-nucleotide polymorphism data". *Genetics* 165 (4): 2213-2233. doi:10.1093/genetics/165.4.2213.
- Ernst, Jason; Kellis, Manolis. 2012. "ChromHMM: automating chromatin-state discovery and characterization". *Nature Methods* 9 (3): 215-216. doi:10.1038/nmeth.1906. PMC 3577932. PMID 22373907.
- Qurbonova, M., Abjalova, M. va boshq. 2021. *O'zbek tili o'zlashma so'zlarining urg'uli lug'ati*. O'quv-uslubiy lug'at. Toshkent: Nodirabegim.
- https://en.wikipedia.org/wiki/Hidden_Markov_model.
- <https://www.freecodecamp.org/news/an-introduction-to-part-of-speech-tagging-and-the-hidden-markov-model-953d45338f24>.
- <https://uzjournals.edu.uz/namdu/vol1/iss6/46>.
- <https://coderlessons.com/tutorials/akademicheskii/obrabotka-estestvennogo-iazyka/pometka-chasti-rechi-pos>.
- <https://habr.com/ru/post/125988>.
- https://ru.wikipedia.org/wiki/Частеречная_разметка.
- https://en.wikipedia.org/wiki/Part-of-speech_tagging#:~:text=In%20corpus%20linguistics%2C%20part%2Dof,its%20definition%20and%20its%20context.
- https://ru.wikipedia.org/wiki/Часть_речи.

The Task of Tagging Part of Speech in Natural Language Processing

Manzura Abjalova¹

Dilrabo Elova²

Abstract

Natural language processing is a complex process that requires consideration of all phenomena, general and particular aspects, exceptions, phonetic, morphological, lexical, grammatical, semantic and even orthoepic features of a social language. At the same time, since the creation of an artificial intelligence system is an important result, it is important to achieve text processing units using digital technologies. In this process, it is necessary to determine the parts of speech of words. In computational linguistics, it is known that this linguistic and technical analysis is a stage of marking words – PoS-tagger, that is, tagging a part of speech of words, and is also considered the main stage of automatic text processing. Initial attempts to create tags in creating a language corpus solved many of the most pressing text problems today. In natural language processing (NLP), word labeling is also a primary concern, which leads to the analysis of various linguistic ambiguities in the text, such as the definition of homonymy and semantics of ambiguous words.

This article discusses about part-of-speech tagging words, its significance in the process of computer analysis of texts, and the methods of tagging.

Key words: *natural language processing, tagging, tag, part-of-speech, formal language, pragmatic feature, corpus, polysemy, homonymy, PoS-tagging.*

References

- Abjalova, M. 2020. *Tahrir va tahlil dasturlarining lingvistik modullari*. Toshkent: Nodirabegim.
- Abjalova, M., Yuldashev, A. 2021. "Methods for Determining Homonyms in Linguistic Systems". *ACADEMICIA: An International Multidisciplinary*

¹Manzura A. Abjalova – Doctor of Philosophy in Philology (PhD), Tashkent State University of Uzbek Language and Literature named after Alisher Navo'i.

E-mail: abjalovamanzura@navoiy-uni.uz

ORCID ID: 0000-0002-1927-2669

²Dilrabo Q. Elova – teacher, Tashkent State University of Uzbek Language and Literature named after Alisher Navo'i.

E-mail: dilrabo@navoiy-uni.uz

ORCID ID: 0000-0002-1927-2968

For citation: Abjalova, A.A., Elova, D.Q. 2021. "The Task of Tagging Part of Speech in Natural Language Processing". *Uzbekistan: Language and Culture* 1: 6–20.

Research Journal 11 (2): 700-715. DOI: 10.5958/2249-7137.2021.00522.X.

- Abjalova, M. 2021. "O'zbek tili Milliy korpusida so'zshakllarni leksikografik baza asosida qidiruv imkoniyatlari". *Kompyuter lingvistikasi: muammo, yechim, istiqbollari*, 12-17. Toshkent: ToshDO'TAU.
- Abjalova, M. 2021. "O'zbek tili milliy korpusida sinonimayzer yoxud sinonimizatorni yaratish masalasi". *O'zbek Milliy va ta'limiy korpuslarining yaratishning nazariy hamda amaliy masalalari*, 38-40. Toshkent: ToshDO'TAU.
- O'zbek tilining izohli lug'ati*. 2006. 80 000 dan ortiq so'z va so'z birikmasi. A. Madvaliev tahriri ostida. 5 jildli. Toshkent.
- Brill, E. 1995. "Transformation-Based Error-Driven Learning and Natural Language Processing: A Case Study in Part-of-Speech Tagging". *Computational Linguistics* 21: 543-565.
- Brill, E. 1992. "A simple rule-based part of speech tagger". *Proceedings of ANLC* 154.
- Baum, L. E.; Sell, G. R. 1968. "Growth transformations for functions on manifolds". *Pacific Journal of Mathematics* 27 (2) 211-227.
- Ganiyeva, Dildora. 2019. "Mazmunij sinkretizm va polifunkcionallik". *NamDU ilmiy axborotnomasi* 6: 275-278.
- Rizayev, S. 2006. *O'zbek tilshunosligida lingvostatistika asoslari*. Toshkent: Fan.
- Starner, Thad, Pentland, Alex. 1995. *Real-Time American Sign Language Visual Recognition From Video Using Hidden Markov Models*. Master's Thesis, MIT, Program in Media Arts.
- Li, N; Stephens, M. 2003. "Modeling linkage disequilibrium and identifying recombination hotspots using single-nucleotide polymorphism data". *Genetics* 165 (4): 2213-2233. doi:10.1093/genetics/165.4.2213.
- Ernst, Jason; Kellis, Manolis. 2012. "ChromHMM: automating chromatin-state discovery and characterization". *Nature Methods* 9 (3): 215-216. doi:10.1038/nmeth.1906. PMC 3577932. PMID 22373907.
- Qurbonova, M., Abjalova, M. va boshq. 2021. *O'zbek tili o'zlashma so'zlarining urg'uli lug'ati*. O'quv-uslubiy lug'at. Toshkent: Nodirabegim.
- https://en.wikipedia.org/wiki/Hidden_Markov_model.
- <https://www.freecodecamp.org/news/an-introduction-to-part-of-speech-tagging-and-the-hidden-markov-model-953d45338f24>.
- <https://uzjournals.edu.uz/namdu/vol1/iss6/46>.
- <https://coderlessons.com/tutorials/akademicheskii/obrabotka-estestvennogo-iazyka/pometka-chasti-rechi-pos>.
- <https://habr.com/ru/post/125988>.
- https://ru.wikipedia.org/wiki/Chasterechnaja_razmetka.
- https://en.wikipedia.org/wiki/Part-of-speech_tagging#:~:text=In%20corpus%20linguistics%2C%20part%2Dof,its%20definition%20and%20its%20context.
- https://ru.wikipedia.org/wiki/Chast'_rechi.