

# ქართული ენის მოდელირების მცდელობა

კოტორეიშვილი ვახტანგი, კურცხალია დაჩი, სტურუა საბა

[vkoto17, dkurt17, sstur17]@freeuni.edu.ge

Free university of Tbilisi

Tbilisi, Georgia

## აბსტრაქტი

ნაშრომში განხილულია ქართული ენის მოდელირების მცდელობა. დავინყეთ მონაცემების შეგროვებით, მათი გაანალიზებითა და ვიზიუალიზაციით. შემდგომ გადავედით დამუშავებულ მონაცემებზე სხვადასხვა ტიპის მოდელის აგებაზე. დავინყეთ **N-Gram**-ით, რომელიც სიტყვათა მიმდევრობების სიხშირებს იყენებს მოდელირებისთვის. მოცემულმა მიდგომამ ვერ შეძლო ქართული ენის სპეციფიკის გადმოცემა. ამის შემდეგ ვცადეთ უფრო კომპლექსური ტრანსფორმერ მოდელები, შევქმენით ვორდ ემბედინგები, გამოვიყენეთ სხვადასხვა ტიპის ტოკენიზაცია. საბოლოოდ კი ბევრი რამ გავიგეთ მოცემული ამოცანისა და მისი სირთულეების შესახებ.

## შესავალი

ჩვენი ამოცანაა ქართული ენის მოდელირება. ენის მოდელის კონცეფცია დიდი ხანია არსებობს კომპიუტერულ მეცნიერებაში. დანყებული სტატისტიკური მოდელებიდან, დამთავრებული ტრანსფორმერებით ბევრნაირადაა შესაძლებელი მოცემული ამოცანის გადაჭრა.

გამოცდილებამ გვაჩვენა, რომ კარგი ენის მოდელის შექმნას ჭირდება ენის დიდი კორპუსი და ასევე დიდი გამოთვლითი რესურსები. მაგალითისთვის, **Open AI**-ს ენის მოდელს **gpt-3**-ს აქვს **175** მილიარდი პარამეტრი. ამ კუთხით, შეიძლება ითქვას, ქართული ენა დაჩაგრულია. ჩვენ **low resource** ენებში გავდივართ და ამას ისიც ემატება, რომ ამ პროექტის განმავლობაში არ გვქონია ნორმალური **GPU** სერვერი (**Kaggle**-ს ჯიფიუს ვიყენებდით)

როდესაც ამოცანაზე ფიქრი დავინყეთ, მივხვდით, რომ რამდენიმე გამოწვევის წინაშე ვიდექით: დავინყოთ იმით, რომ გამოთვლითი რესურსების სიმცირის გამო, მოგვინია მცირე ზომის კორპუსის გამოყენება; ამასთანავე, რადგან ქართულ ენაში არ გაგვაჩნია ხელსაწყოები, რომლითაც შევძლებთ სიტყვის შედარებით მარტივ ფორმაზე დაყვანას (მაგალითად ფუძის გამოყოფა რომ შეიძლებოდეს) შედეგად მივიღეთ ძალიან ბევრი უნიკალური ტოკენი, რომელთა სიხშირეც საკმაოდ მცირე იყო; ამას ემატებოდა ისიც, რომ თითქოს ყველაზე მიმზიდველი კორპუსი - ვიკიპედია - შეიცავს სპეციფიკურ ტექსტებს და, ჩვენი აზრით, კარგად არ ასახავს ქართულ ენას (ტოპ 10 ყველაზე ხშირ სიტყვაში შედის “ნელი” და “წლის”);

მიუხედავად ზემოხსენებული გამოწვევებისა, ვცადეთ ტრივიალური მიდგომით დაგვეწყო და ნელ-ნელა “გაგვერთულებინა”, პროცესში გვესწავლა რას ვაკეთებდით არასწორად და როგორ შეიძლებოდა მიღებული შედეგის გაუმჯობესება. ამ ყველაფერზე უფრო ვრცლად მომდევნო თავებში დავწერთ.

## მონაცემების შეგროვება

როგორც ვახსენეთ, ენის მოდელს დიდი რაოდენობით ტექსტური მონაცემი ესაჭიროება. ვცადეთ დაკავშირება ქართული ენის ეროვნული კორპუსის [საიტთან](#), თუმცა პასუხი არ გაგვცეს. კორპუსის პოვნის მცდელობა ნაკლებად რელევანტური გახდა მის მერე, რაც გავიაზრეთ, რომ ნორმალურ GPU სერვერს ვერ ვიპოვიდით. საბოლოოდ, სხვა გუნდებთან ერთად შევჯერდით ვიკიპედიას კორპუსზე, რომლის მოპოვებაც საკმაოდ მარტივი გამოდგა. თუმცა, მოდელის გასატესტად გავპარსეთ [marketer.ge](#).

## მონაცემთა ანალიზი და პრეპროცესინგი

როგორც შესავალში აღვნიშნეთ, კორპუსი არ შეესაბამება ქართული ენის ზოგად სახეს. ვიკიპედიას სტატიების სპეფიციკიდან გამომდინარე, გვხვდება მეტი სასვენი ნიშანი. ამავდროულად, 10 ყველაზე ხშირ სიტყვაში შედის ისეთი სიტყვები, როგორიცაა “წელი” და “წლის”, რაც ნამდვილად არ არის კორექტული ქართული ენისთვის.

ტექსტის პრეპროცესინგის პირველი და ტრივიალური მიდგომა იყო უბრალოდ სასვენი ნიშნების გამოყოფა და “სვენიების” მიხედვით გაყოფა. ამ მიდგომით დავთვალეთ უნიკალური ტოკენები, ავარჩიეთ ტოპ 80-90 ათასი, რომელიც შევიყვანეთ ჩვენს **vocabulary**-ში და დანარჩენს მივანიჭეთ **unk** ტოკენი.

მოცემული პრეპროცესინგის შედეგად მივიღეთ ძალიან ბევრი **unk** ტოკენი, რაც, თავისთავად, პრობლემაა. სამწუხაროდ, ეს პრობლემა ვერ მოგვარდებოდა **vocabulary**-ს ზომის გაზრდით, რადგან კორპუსი მცირე ზომის გვექონდა, **vocabulary**-ს გაზრდა კი ავტომატურად ნიშნავდა მოდელის პარამეტრების ზრდას და ისეთი სიტყვების გათვალისწინებას, რომლებიც იშვიათად გვხვდება ისედაც მცირე ზომის კორპუსში. ჩვენი აზრით, **vocabulary**-ს გაზრდა ვერ მოაგვარებდა ზემოხსენებულ პრობლემას. შესაბამისად, გადავწყვიტეთ ყურადღება გადაგვეტანა პრეპროცესინგის ნაწილზე.

პრეპროცესინგის მეორე მცდელობისას ვეცადეთ **unk** ტოკენში უფრო მეტი ინფორმაცია ჩაგვედო. პირველი მცდელობა ერთსა და იმავე შემთხვევად განიხილავდა უცხოურ სიტყვას, ქართულ უცნობ სიტყვას, რიცხვს და ა.შ. შევქმენით სხვადასხვა ტიპის **unk** ტოკენები : ციფრები (**NUM\_SML**), რიცხვები 10დან 2100მდე (**NUM\_MED**), რიცხვები 2100დან (**NUM\_LRG**), უცნობი სიმბოლო (**SYMB\_UNK**), უცნობი უცხოური სიტყვა (**FOREIGN\_UNK**), უცნობი ქართული სიტყვა (**FOREIGN\_UNK**).

გარდა ნახსენები პრეპროცესინგისა, ვცადეთ **subword** ტოკენაიზერის (**BPE Tokenizer**) გამოყენება, რამაც ბევრად უკეთესი შედეგი მოგვცა. ეს ბუნებრივია, რადგან აქამდე თუ ორი მსგავსი სიტყვა, მაგალითად “ადამიანი” და “ადამიანები” ორ სხვადასხვა ტოკენს წარმოადგენდა, **subword** მოდელირებით შესაძლებელია ამ ორ სიტყვას საერთო ტოკენი/ტოკენები ჰქონდეს.

## მოდელის აგება

სანამ ძირითადი მოდელის აგებას დავიწყებდით, ავაგეთ მარტივი **N-gram** მოდელი. **N-gram**<sup>1</sup> მოდელი გულისხმობს ტექსტის სტატისტიკურ ანალიზს შემდეგნაირად: ვითვლით ერთად შემხვედრი სიტყვათა ჯგუფების რაოდენობას. შემდეგ ალბათობების შენონისთვის გამოვიყენეთ **Laplace Smoothing**. იქიდან გამომდინარე რომ **N-gram** მოდელი გულისხმობს ყველა  $n$ -ჯგუფი ტოკენების შენახვას, პრობლემა შეგვექმნა მეხსიერების მხრივ. ამიტომ მოდელი დავატრენინგეთ მონაცემთა მესამედზე. როგორც მოსალოდნელი იყო, “მოდელის” მიერ შექმნილი წინადადებები სატრენინგო დატაზე იყო მთლიანად მორგებული და თავისით ვერ იგონებდა ახალ, გამართულ წინადადებებს. ამას დავმატა ისიც, რომ მცირე ზომის კორპუსის გამო, მოგვინია **N** მცირე აგველო და, შესაბამისად, დაგენერირებულ ტექსტს არ გააჩნდა კონტექსტი.

გადავწყვიტეთ, უფრო კომპლექსური მიდგომა გვეცადა. ჩვენთვის ძალიან საინტერესო იყო ტრანსფორმერ არქიტექტურა, რაც დღეს **SOTA**-დ ითვლება. შესაბამისად, დავიწყეთ ყველაზე მარტივი ვერსიით, ანუ არქიტექტურით, რომელიც ნაშრომში - **Attention is all you need** - იყო შემოთავაზებული.<sup>2</sup>

ტექსტის გენერაციისთვის გამოვიყენეთ **Beam Search** ალგორითმი. ბიმ-სერჩის დროს ტექსტის გენერაციის დროს ვაკვირდით **N** ყველაზე მაღალქულიან შედეგს, რომელსაც შემდეგი სიტყვის გენერაციის დროს ვიყენებთ. ყოველ ნაბიჯზე “ვსხლავთ” და ვტოვებთ მხოლოდ **N** მაღალქულიან შედეგებს. ბიმ სერჩის გარდა გამოვიყენეთ **greedy** და რენდომ გზებიც თუმცა ვფიქრობთ ბიმ სერჩი ყველაზე გამოსადეგია.

გადავწყვიტეთ მოდელს “დავხმარებოდით” და შევქმენით ვორდ ემბედინგები **word2vec** ალგორითმის მიხედვით. ქვემოთ მოყვანილი მოდელები დავატრენინგეთ როგორც **pretrained** ემბედინგებით, ასევე მათ გარეშე, თუმცა რეალური სხვაობა შედეგში არ ქონია. ემბედინგებზე მეტი ინფორმაციისთვის წაიკითხეთ ნოუტბუქი - **Word2vec Visualization and Analysis.ipynb**.

პირველი მოდელი დავატრენინგეთ პირველი პრეპროცესინგით დამუშავებულ ტექსტზე შემდეგი პარამეტრებით: ემბედინგის განზომილება - **100**; ენკოდერში ჰიდენ ლეიერის ზომა - **200**; ენკოდერ ლეიერების რაოდენობა - **2**; **attention**-ების რაოდენობა - **2**; დროფაუთი - **0.3**; **Cross Entropy** ობჯექტივი ფუნქცია და **SGD** ოპტიმაიზერი. შედეგად მივიღეთ მოდელი, რომელიც ძალიან დიდ უპირატესობას ანიჭებდა **unk** ტოკენს. ეს ბუნებრივიც იყო გამოყენებული პრეპროცესინგიდან გამომდინარე.

როდესაც იგივე მოდელი მეორე ტიპის პრეპროცესინგით დამუშავებულ ტექსტზე ვატრენინგეთ, უკეთესი, მაგრამ მაინც ცუდი შედეგი მივიღეთ. მოდელი ისევ უპირატესობას ანიჭებდა **unk** ტოკენებსა და ამჯერად უკვე სასვენ ნიშნებსაც. ამის მერე, ვივარაუდეთ, რომ პრობლემა კორპუსში იყო: გვხვდებოდა დიდი რაოდენობით სასვენი ნიშნები და სიტყვების უმეტესობა მცირე რაოდენობით იყო წარმოდგენილი.

<sup>1</sup> Daniel Jurafsky & James H. Marti. Speech and Language Processing.- Chapter 3- N-gram models

<sup>2</sup> Ashish Vaswani Noam Shazeer Niki Parmar Jakob Uszkoreit Llion Jones Aidan N. Gomez Lukasz Kaiser Illia Polosukhin - Attention is All You Need

ზემოხსენებული პრობლემის გადაჭრა ვცადეთ ოპტიმიზერის შეცვლით. SGD-ს მაგივრად, არჩევანი შევაჩერეთ **Adagrad**-ზე<sup>3</sup>. SGD-სგან განსხვავებით, **Adagrad**-ი განსხვავებულ ლერნინგ რეითს ურჩევს თითოეულ პარამეტრს თითოეულ თრენინგ ნაბიჯზე. ეს კი იმაში გვეხმარება, რომ ხშირი სიტყვები დაბალი ლერნინგ რეითით შეიცვლებიან, ხოლო იშვიათები - დიდით. სამწუხაროდ, არც ეს ცვლილება გამოდგა მნიშვნელოვანი და მიღებულ მოდელს ისევ უჭირდა ადეკვატური ტექსტის დაგენერირება.

გამომდინარე იქედან, რომ გვექონდა მცირე გამოთვლითი და დროითი რესურსები, გადაწყვიტეთ შეგვეცვალა მიდგომა. **Huggingface** ბიბლიოთეკის გამოყენებით შევქმენით **subword** ტოკენაიზერი და **Roberta**-ს არქიტექტურის მოდელი ვატრენინგეთ ვიკიპედიას კორპუსზე<sup>4</sup>. ამ მიდგომამ ბევრად უკეთესი შედეგი მოგვცა, თუმცა იგი მაინც შორსაა სასურველისგან. მოდელი მაინც უპირატესობას ანიჭებს სასვენ ნიშნებს და დაგენერირებული ტექსტი კარგად არ ასახავს ქართულ ენას.

## ევალუაცია

ევალუაციისთვის გამოვიყენეთ ქართული ვებსაიტის - **Marketer.ge**-დან წამოღებული ქართული სტატიები. მარკეტერის სტატიები ვიკიპედიის სტატიებისგან განსხვავებული ლექსიკით არის დაწერილი. შესაბამისად ევალუაციისას მიღებული განსხვავებები და შედეგი კარგი დასაკვირვებელია. იქიდან გამომდინარე, რომ მოდელის ტრენინგმა პრობლემები შეგვიქმნა, საშუალება არ მოგვეცა რომ ევალუაციისას კარგად შეგვესწავლა მოდელის ქცევა.

## დასკვნა

მოცემული კვლევით ფეხი შევდგით იმ გაუკვალავ ტყეში, რომელსაც ქვია ქართული ენის მოდელირება. ამოცანა, მართლაც რომ, რთულია, თუმცა ვფიქრობთ, რომ დღევანდელი NLP-ს ტექნიკებით სავსებით შესაძლებელია კარგი ენის მოდელის შექმნა.

დავინწყეთ ტრივიალური მიდგომით, ვცადეთ დაგვენახა პრობლემები, შეგვეცვალა მიდგომა და ასე გაგვეუმჯობესებინა შედეგი, მაგრამ სასურველი შედეგის მიღება ჯერჯერობით ვერ შევძელით.

ნათლად დავინახეთ ის, რომ კარგი ენის მოდელის შექმნას ესაჭიროება დიდი რაოდენობით მონაცემები და გამოთვლითი რესურსი. ჩვენი აზრით, სწორ გზაზე ვართ და თუ მეტ ყურადღებას დავუთმობთ ენის კარგი კორპუსის შექმნას, საკმაოდ კარგი შედეგი გვექნება.

## ბიბლიოგრაფია

<sup>3</sup> [An overview of gradient descent optimization algorithms](#)

<sup>4</sup> Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, Veselin Stoyanov - RoBERTa: A Robustly Optimized BERT Pretraining Approach

- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, Veselin Stoyanov - RoBERTa: A Robustly Optimized BERT Pretraining Approach
- Ashish Vaswani Noam Shazeer Niki Parmar Jakob Uszkoreit Llion Jones Aidan N. Gomez Lukasz Kaiser Illia Polosukhin - Attention is All You Need
- Daniel Jurafsky & James H. Marti. Speech and Language Processing.- Chapter 3- N-gram models
- [An overview of gradient descent optimization algorithms](#)