

# Muturretik muturrerako solasaldi sistema

Maitane Martinez Eguluz

## Abstract

Proiektu honetan, dialogoa mantentzeko muturretik muturrerako bi sistema aurkezten dira. Lehenengoa, ingelesez garatu da, bigarrena, berriz, euskaraz. Ingeleseko sistema Telegram Bot batera egokitu da.

## 1 Sarrera

Proiektu honen ataza elkarrizketa mantentzen duen sistema bat sortzea da. Horretako, muturretik muturrerako sare errekurrentea da erabiltzen den sistema. Entrenatzeko datuak filmetako azpтитuluak dira. Bi hizkuntzetan garatu da sistema: ingelesez eta euskaraz. Gainera, Telegramen Bot bezala egokitu da ingelesez entrenatua izan den sistema.<sup>1</sup>

## 2 Erlazionatutako lanak

Banadanu-k eta beste batzuek neurona sareetan oinarritutako itzulpen automatikoa proposatu zuten. Itzulpen tradizionalak ez bezala, kodetzaile-deskodemzaile motako modeloa proposatzen dute ingelesetik frantseserako itzulpena egiteko. Kodetzaile sare neuronalak iturriko esaldia irakurri eta luzera finkoko bektore batean kodetzen du. Dekodemzaileak kodetutako bektorearen itzulpena egiten du. (Bahdanau, 2014)

Lison-ek eta beste batzuek, artikulu horretan, azpтитuluen aurreprozesazioan eta lerrokatzean hobekuntza ugari proposatzen dituzte, hala nola OCR akatsak automatikoki zuzentzea. Gainera, metadatuak erabiltzen dituzte azpтитuluen kalitatea estimatzeko eta azpтитulu bikoteak ebaluatzeko. (Lison, 2016)

## 3 Sistema

Atal honetan muturretik muturrerako solasaldi sistema azaltzen da. Sistema hau Bahdanau-k eta beste batzuek proposatutako lanean oinarritzen da (Bahdanau, 2014). Horretarako, dialogoa itzulpen ataza bat bezala proposatzen da. Esate baterako:

- Itzulpen automatikoa :
  - Sistemaren sarrera: Egun on guztioi
  - Sistemaren irteera: Buenos días a todos
- Dialogoa:
  - Sistemaren sarrera: Egun on guztioi
  - Sistemaren irteera: Baita zuri ere

Sistema kodetzaile-deskodemzaile motako modeloa da. Kodetzaileak iturriko esaldia irakurtzen du bi norabideko sare errekurrentearen bidez eta luzera finkoko bektore batean kodetzen du. Deskodemzaileak, kodetutako bektorea erabiliz, esaldi berria sortzen du.

### 3.1 Kodetzailea

Kodetzaileak, sarrerako esaldia irakurtzen du tokenen bektorea bezala,  $x = (x_1, \dots, x_{T_x})$ . Horretako, bi norabideko sare errekurrentea erabiltzen da. BiRNN sarea aurreranzko eta atzerako RNN-z osatutako dago. Aurreranzko RNN-a  $\vec{f}$ , ezkererik eskuinera irakurtzen du esaldia eta aurreranzko geruza ezkutua kalkulatu du  $(\vec{h}_1, \dots, \vec{h}_{T_x})$ . Atzerazko RNN-ak  $\overleftarrow{f}$ , eskuinetik ezkerreko irakurtzen du esaldia eta atzerazko geruza ezkutua kalkulatu du  $(\overleftarrow{h}_1, \dots, \overleftarrow{h}_{T_x})$ .

Hitz bakoitzeko ohar bat lortzen da aurrerazko geruza eskutua  $\vec{h}_j$  eta atzeranzko geruza ezkutua  $\overleftarrow{h}_j$  konkatatuz:  $h_j = [\vec{h}_j^T; \overleftarrow{h}_j^T]^T$ . Horrela, oharak  $h_j$  aurreko eta ondorengo hitzen informazioa edukiko du. Oharrek testuinguru bektorea osatzeko balio dute.

<sup>1</sup>Lanaren GitHub helbidea: <https://github.com/MaitaneMartinez/DIAL->

### 3.2 Atentzio geruza

Testuinguru bektorea  $c_i$  oharren  $h_i$  batuketaren bidez kalkulatzen da:

$$C_i = \sum_{j=1}^{T_x} \alpha_{ij} h_j. \quad (1)$$

ohar bakoitzaren  $h_j$  pisua  $\alpha_{ij}$  honela kalkulatzen da:

$$\alpha_{ij} = \frac{\exp(e_{ij})}{\sum_{k=1}^{T_x} \exp(e_{ik})}, \quad (2)$$

non

$$e_{ij} = a(s_{i-1}, h_j) \quad (3)$$

$s_{i-1}$  deskodetzailearen RNN eskutuzko geruza da (5) eta  $j$  sarrerako esaldiaren oharra  $h_j$  da. a lerrotatze-eredua entrenatua izan da beste os-agaiekin batera.

### 3.3 Deskodetzailea

Deskodetzailea, testuinguru bektorea  $c$  eta aurretik iragarritako hitza  $y_{i-1}$  emanik, hurrengo hitza  $y_i$  iragartzeko entrenatzen da. Horretarako, probabilitate bektorea sortzen da:

$$p(y_i | y_1, \dots, y_{i-1}) = g(y_{i-1}, s_i, c_i) \quad (4)$$

non  $s_i$  RNN-ko eskutuzko geruza den, i denboran:

$$s_i = f(s_{i-1}, y_{i-1}, c_i) \quad (5)$$

$f$  eta  $g$  funtzio ez linealak dira. Hiru dekodetze estrategia desberdin erabiltzen dira: top1, topk eta multinomial. top1 estrategiak ikasitako token posibleena itzultzen du. topk estrategiak k erantzun posibleenetatik bat itzultzen du. Sistema honetan  $k=3$  definituta dago. Azkenik, multinomial estrategia, banatze nominala erabiltzen du tokenak aukeratzeko.

Telegram Bot-a sortu da GitHub dauden jarraibideetan oinarrituta <sup>2</sup>. Ingeleseko sistema erabiltzen da eta aldagai global baten bitartez deskodetze estrategia alda daiteke, aldagai honi balio desberdinak esleitzuz.

<sup>2</sup>Your first Bot: <https://github.com/python-telegram-bot/python-telegram-bot/wiki/Extensions-%E2%80%93-Your-first-Bot>

## 4 Datuak

Atal honetan, erabilitako corpusak eta datuak azaltzen dira. OpenSubtitles corpora erabiltzen da bi sistemak sortzeko <sup>3</sup>. Corpus honetan 92 hizkuntzako azpitituluak daude. Gure kasuan, ingelesezko eta euskarazko azpitituluak erabili dira soilik.

Lehenengo sisteman, ingelesezko filmen esaldiak erabiltzen dira. Esaldiak binaka elkartzen dira tabulazio baten bidez banatuta: bigarren esaldia lehenengoaren erantzuna litzateke. Guztira, milio bat esaldi pare daude. Baina, batzuetan galdera eta erantzuna ez datoz bat, adibidez:

- Galdera: Friends.
- Erantzuna: Good.

Bigarren sisteman <sup>4</sup>, euskarazko filmen esaldiak erabiltzen dira. Esaldi horiek binaka jartzen dira, ingelesezko egiturari jarraituz. Gainera, garbiketa bat aplikatzen da esaldietan. Garbiketari, dialogo "-” ikurra eta puntu finala ”.” kentzen dira.

## 5 Emaizak

Atal honetan, esperimentuak eta horien emaitzak azaltzen dira. Horretarako, ingeleserako entrenatu izan den muturretik muturrerako solasaldi sistema probatzen da. Hiru frogak desberdin egin dira, dekodetze estrategia desberdinak erabiliz.

Elkarrizketa sinple bat erabili da probak egiteko:

- Hello
- How are you?
- So, what are your plans for this weekend?
- Bye

Lehenengo sarreratik hirugarrenera zailtasuna handitzen dira. Azken sarrera erabiltzen da sistema amaitzeko. Honako erantzun desberdinak hauek lortu dira:

<sup>3</sup>OpenSubtitles web gunea: <http://opus.nlpl.eu/OpenSubtitles-v2018.php>

<sup>4</sup>Sistema deskargatzeko helbidea: <https://drive.google.com/drive/folders/1i09cLB1rsKBMUbpRQzJTbw2J308WjaEW?usp=sharing>

### top1

- Hello ,
- I ' m fine .
- I ' m not .
- Bye ,

Beti, erantzun berdinak itzultzen dira. Erantzunak laburrak dira eta, hirugarren galderan, berriz, sistemak ez du sortzen erantzun egokirik. Gainera, puntuazio marka ez zuzenak erabiltzen ditu.

### topk

- I , hello .
- Fine , fine , thanks , man ! .
- I ' s . to go to bed ? ?
- I , you

Erantzun desberdinak itzultzen ditu. Aurreko estrategiarekin konparatuz, erantzunak luzeagoak dira, baina zentzu gabeko puntuazio markak itzultzen ditu.

### multinomial

- Hello you .
- I ' m fine
- Yes .
- Good bye

top1 estrategiaren antzeko erantzunak itzultzen ditu. Kasu honetan, puntuazio markak egokiagoak dira.

Ingelesezt erabilitako elkarrizketa euskaratu da bigarren sistema testatzeko:

- Kaixo
- Zer moduz zaude?
- Zer egingo duzu asteburuan?
- Agur

Honako erantzun hauek bueltatzen ditu sistemak:

### top1

- Zer moduz?
- Ondo
- < unk >
- Agur

Hirugarren galderarentzat ez dauka erantzun posiblerik.

### topk

- Kaixo,
- Oso ondo,
- Ez dakit zer
- Agur esan

Erantzun posibleenetatik aiposenak eta top1 erantzunekin desberdinak direnak aukeratu dira.

### multinomial

- Kaixo
- Oso ongi
- Ez dut
- Agur

topk bezala, erantzun guztietatik aiposenak aukeratu dira. Hiruretan antzeko erantzunak itzultzen ditu sistemak.

Azkenik, 1. irudian ikusi daiteke Telegram Bot-a erabiliz mantendu den elkarrizketa.



Figure 1: Telegram Bot-arekin izandako elkarrizketa

/star bidez dialogoa hasiko da eta dekodetze estrategia desberdinak erabili daitezke elkarrizketan zehar "/" aurretik jarrita. Ingeleseztako sistemak dituen arazo berdinak edukiko ditu Telegrameko Bot-ek.

## 6 Analisia

Lehenengo sisteman lortutako emaitzak ikusita, esan daiteke hirugarren galderaren erantzunak ez direla zentzuzkoak. Hau da, sistema ez da gai galdera honi erantzuna emateko. Gainera, multinomial estrategian izan ezik, beste bietan puntuazio marka desegokiak erabiltzen ditu. Hori gertatzen da garbiketa sakon bat falta delako puntuazio marketan. Gainera, esaldi laburrak itzultzen ditu sistemak; izan ere, token ohikoenek pisu handia edukiko dute eta horiek bakarrik itzultzeko joera edukiko du sistemak. Horien artean, puntuazio markak daude.

Bigarren sisteman, puntuazio marken arazoak konpondu egin dira erantzun gehienetan. Gainera, hirugarren galderarentzat topk estrategia erabiliz erantzun zuzen bat itzultzen du. Baina, gutxi agertzen diren hitzak ez ditu ikasten, adibidez, izen propioak. Horren ordez, *< unk >* motako tokenak itzultzen ditu. Aproposa izango litzateke euskaraz entrenatutako tokenizatzailerak erabiltzea, baita normalizatzailerak sistema hobeto funtzionatzeko.

Sistemak ez du dialogo bat mantentzen, galde-tutakoari baino ez dio erantzuten, aurreko elkar-rizketa kontuan hartu gabe.

## 7 Ondorioak

Proiektu honetan, bi sistema aurkezten dira. Biak muturretik muturrerako solasaldi sistemak dira, baina lehenengoa ingelesezko filmen azpitituluekin entrenatua eta bigarrena euskarazkoekin. Gainera, Telegram Bot-a sortu da ingelezko sistema erabiliz. Sistema horiek hiru deskodetze estrategia erabiltzen dituzte: top1, topk eta multinomiala.

Etorkizuneko lanei begira, interesgarria litzateke euskarazko sistema entrenatzea elkarriketen bidez bakarrik. Datu horietan, galdera-erantzun motako egiturak errespetatzen dira. Beraz, nahiz eta datu kopurua txikiagoa izan, gerta daiteke erantzun aproposagoak itzultzea.

Bi hizkuntzako sistema sortzea izango litzateke beste proposamen bat. Entrenatzeko datuetan errespetatutako beharko litzateke galderaren berdina izatea erantzunaren hizkuntza. Beraz, testatzeko orduan, sistemari hitz egiten diozun hizkuntza errespetatu beharko luke eta ez nahastu.

## References

- Cho K. Bengio Y. Bahdanau, D. 2014. Neural machine translation by jointly learning to align and translate. arXiv:1409.0473.
- Tiedemann J. Lison, P. 2016. Opensubtitles2016: Extracting large parallel corpora from movie and tv subtitles.