



Universitatea Transilvania din Braşov
Facultatea de Matematică şi Informatică
Specializarea Tehnologii Moderne
în Ingineria Sistemelor Soft

LUCRARE DE DISERTAȚIE

Autor: Muşat Bogdan-Adrian
Coordonator: Conf. dr. Lucian-Mircea Sasu

Braşov
2017

Cuprins

1	Introducere	2
1.1	Motivația alegerii temei	4
1.2	Structura lucrării	5
1.3	Acknowledgement	5
2	Lucrări similare	6
2.1	Modele de regăsire	6
2.1.1	Duolingo	7
2.1.2	Facebook Messenger	7
2.2	Modele generative	8
3	Arhitectura	9
3.1	Rețele neurale artificiale	9
3.2	Rețele neurale artificiale recurente	10
3.3	Modele Sequence-to-Sequence	12
3.4	Hierarchical Recurrent Encoder-Decoder	13
3.4.1	Celula GRU	14
3.4.2	HRED bidirecțional	15
3.5	Word Embeddings	15
3.6	Predicție	16
3.6.1	Greedy	16
3.6.2	Beam search	16

Capitolul 1

Introducere

În zilele noastre, societatea se confruntă tot mai mult cu diverse probleme tehnice în mediul online, industrial, cotidian etc. Până de curând, aceste probleme puteau fi rezolvate doar cu ajutorul unui specialist în domeniul respectiv. După cum știm, omul reprezintă o resursă limitată când vine vorba de suport tehnic, deoarece cererile pot fi numeroase, în consecință răspunsul primit poate fi întârziat. De asemenea, trebuie luat în considerare costul prohibitiv de întreținere a unor persoane responsabile cu acest tip de suport.

Recent, diverse companii de succes precum Google [1], Facebook[2], Microsoft [3], Apple [4] și alții au făcut un demers spre adoptarea unor sisteme inteligente de comunicație numite chatbots. În cazul chatbots-ilor orientați pe suport tehnic, scopul lor este tocmai de a imita sprijinul pe care o persoană reală îl poate oferi. Un sistem computațional inteligent poate fi folosit pentru a răspunde multor cereri simultane, iar costurile de întreținere sunt mici practic, odată ce sistemul este dezvoltat, este necesară doar expunerea lui, de exemplu ca serviciu web. Cu cât tehnologia și munca de cercetare în această arie progresează, cu atât sistemele de suport automat devin din ce în ce mai inteligente și mai apropiate de ce un om este capabil să ofere [19-21]. Se preconizează că în următorii zece ani, aceste sisteme vor fi capabile să înlocuiască cu succes o mulțime de sarcini. Tot așa cum automatizarea industrială de producere a vehiculelor reprezintă un standard în era contemporană, așa vor reprezenta și acești asistenți conversaționali un standard în anii ce vor urma.

Majoritatea modelelor actuale se bazează pe un tip de oferire a unor răspunsuri predefinite. Aceste tipuri de sisteme pot oferi doar soluții existente, nefiind capabile de a genera nimic nou. Pe de altă parte, chatbots-ii

bazați pe inteligență artificială reprezintă cea mai scalabilă modalitate de comunicare între clienți și mediile de afaceri care le sunt dedicate [5].

Vorbim deci de evoluția unor modele de chatbots bazate pe pattern matching către unele bazate pe modele generative, care, precum omul, pot emite răspunsuri bazate pe experiențe anterioare. În ultimii ani, ramura inteligenței artificiale numită machine learning (învățare automată) a luat amploare în această direcție prin curentul numit deep learning [6-8]. Acest curent a produs o mulțime de rezultate spectaculoase atât în direcția procesării naturale de limbaj cât și a procesării de imagini.

Problema asistenților conversaționali este una de procesare a limbajului natural. Sistemul trebuie să fie capabil să "înțeleagă" informația primită de la o persoană și să producă un răspuns cât mai coerent și util. Însă cum înțelege un calculator o limbă? Pentru a răspunde la această întrebare, vom face o analogie la cum învață un om o limbă. Pornește de la anumite cuvinte de bază iar apoi pe baza acestora, învață cuvinte tot mai complexe. Începe să creeze fraze prin care leagă aceste cuvinte precum și gramatica specifică limbajului. Practic, totul se bazează pe o anumită experiență trecută. Conversația este o modalitate ce facilitează și impulsionează deprinderea limbajului: omul este pus în fața unor contexte de utilizare, fenomen ce întărește deprinderea de utilizare a cuvintelor sau expresiilor individuale. Fiecare cuvânt nou învățat reprezintă o experiență spre învățarea în continuare a limbajului. Încă un lucru la care omul este bun este capacitatea de a întreține o conversație lungă cu o altă persoană și să rețină tot felul de informații noi pe parcurs.

În abordările de chatbot actuale se dorește imitarea acestei modalități de învățare pentru un sistem computațional. Modelarea cea mai puternică în ziua de azi pentru o astfel de problemă este oferită de deep learning [6-8]. Folosind un număr mai mare de neuroni față de arhitecturile shallow folosite până la începutul anilor 2000 și plecând de la seturi de instruire masive - în cazul de față corpusuri de text cu cât mai multe fraze în limba pe care o dorim a fi învățată - se formează modele generative care pot produce continuarea unor propoziții, fraze etc. Ideea din spatele abordării deep learning este că sistemul devine mai inteligent pe măsură ce avem tot mai multe date - în acest caz, spețe de conversații. Pentru a procesa o asemenea cantitate de date de instruire este nevoie de putere computațională pe măsură. Asta presupune pe scurt, capacitate hardware. Calculele matematice efectuate pentru învățarea unui limbaj de către un sistem sunt complexe și numeroase (de ordinul sutelor de milioane). Folosirea unui procesor, chiar multicore de

ultimă generație, deși o idee posibilă, este considerată a duce la sugrumarea procesului de instruire.

În locul microprocesorarelor se preferă folosirea de plăci grafice (Graphical Processing Units, GPU). Inițial acestea au fost dezvoltate pentru rulare rapidă de jocuri, dar potențialul lor a fost rapid intuit și exploatat prin programare paralelă. Deoarece placa video conține mult mai multe nuclee de procesare (câteva mii, comparate cu cele 4-8 nuclee tradiționale dintr-un microprocesor actual), este preferată programarea și rularea modelelor computaționale pe GPU. Sunt dezvoltate biblioteci care facilitează unui programator dezvoltarea de aplicații de machine learning pe GPU: Tensorflow [9], Theano [10], Caffe [11], Keras [12] etc.

1.1 Motivația alegerii temei

O potențială utilizare este augmentarea sistemului de creare a tichetelor: Universitatea Transilvania folosește în mod curent un sistem de suport tehnic bazat pe tichete, care apoi sunt procesate de persoane reale pentru rezolvare (Figura 1.1, Figura 1.2). Un tichet presupune primirea de la solicitant a cât mai multor detalii legate de problemă. De regulă, din lipsă de experiență, detaliile furnizate sunt insuficiente pentru o soluționare eficientă, motiv pentru care, după completarea inițială a tichetului se poartă un dialog între suportul IT și solicitant pentru aflarea de informații suplimentare legate de problema raportată. Acest lucru consumă timp; se poate îmbunătăți procesul prin demararea cât mai rapidă a unui dialog solicitant - chatbot prin care să se completeze tichetul cât mai fidel.



Figura 1.1: Sistem pentru suport IT pe portalul Universității



Figura 1.2: Pagina de sesizări pentru biroul IT

Un sistem artificial de chatbot, pe baza dialogurilor anterioare înregistrate și a unei similitudini a cererilor, ar putea fie să solicite mai multe detalii, fie să sugereze pași de rezolvare. În cazul în care există un corpus de cunoștințe (knowledge-base) dat de experiențele anterioare (dialoguri, soluții date), e posibil ca el să fie exploatat în mod automat [19-21].

1.2 Structura lucrării

TODO

1.3 Acknowledgement

TODO

Capitolul 2

Lucrări similare

În anul 1950, matematicianul Alan Turing a propus un test care pune la încercare abilitatea unei mașini computaționale de a manifesta inteligență echivalentă cu cea umană. Un evaluator uman judecă o conversație între un om și o mașină desemnată să genereze răspunsuri cât mai naturale. Evaluatorul este conștient că unul dintre partenerii angrenați în discuție este o mașină. Conversația este limitată doar la text astfel încât rezultatul să nu depindă de abilitatea calculatorului de a genera sunete. Dacă evaluatorul nu poate diferenția mașina de om, putem spune că aceasta a trecut testul. Până în prezent, acest test rămâne încă în picioare.

Există două mari abordări la ora actuală de a genera limbaj cât mai natural. Avem în primul rând modele de tip regăsire, unde răspunsul este stocat într-o sursă de date și returnat pe baza unor metode de pattern matching. Cea de-a doua variantă o reprezintă modelele generative, care produc un răspuns dinamic folosind diverse metode din teoria probabilităților.

2.1 Modele de regăsire

Modelele de acest tip folosesc o sursă de date care conține numeroase răspunsuri predefinite. Răspunsul este ales folosind o metodă euristică pentru o potrivire cât mai bună cu putință, luând în considerare intrarea și contextul. Tipul de euristică folosit poate fi ceva simplu precum o expresie bazată pe o regulă de potrivire sau ceva mai complex cum ar fi un clasificator de Machine Learning. Se poate deduce foarte ușor că aceste sisteme nu generează text nou. O problemă uriașă a acestor modele o reprezintă incapacitatea de a reacționa la cazuri nemaiîntalnite pentru care nu există un răspuns potrivit. Acestea au totuși avantajele lor. Datorită sursei de date cu răspunsuri create

de oameni, aceste metode nu produc erori gramaticale. În prezent, acestea reprezintă abordarea sigură pentru problemele în care răspunsul este unul sensibil, într-un domeniu precum cel medical, de exemplu.

2.1.1 Duolingo

Populara aplicație de învățare a limbilor străine Duolingo (Figura 2.1) folosește o abordare interesantă cu privire la chatbots. Aceasta dorește să își ajute utilizatorii să practice o nouă limbă prin conversații cu chatbots. Având în vedere că o conversație este considerată a fi printre cele mai bune moduri de a învăța o limbă străină, utilizatorii Duolingo pot vorbi cu bot-ul oricât de mult își doresc, iar acesta îi va corecta și le va propune răspunsuri potrivite. Mai mult de atât, poate estima progresul utilizatorului pentru a-și crește nivelul de dificultate, păstrând astfel constantă provocarea.

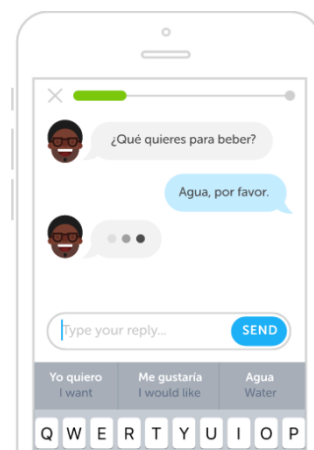


Figura 2.1: Conversație cu chatbot folosind Duolingo

2.1.2 Facebook Messenger

Facebook s-a alăturat întru totul afacerii conversaționale, astfel încât și-a transformat aplicația Messenger într-un business de mesagerie. Compania a integrat plățile peer-to-peer în Messenger în anul 2015, apoi urmând să lanseze un API pentru chatbots, astfel încât business-urile să poată crea interacțiuni pentru clienți. Poți comanda flori, să navighezi printre ultimele trend-uri în materie de modă, să comanzi Uber, toate dinăuntrul chat-ului de Messenger (Figura 2.2).

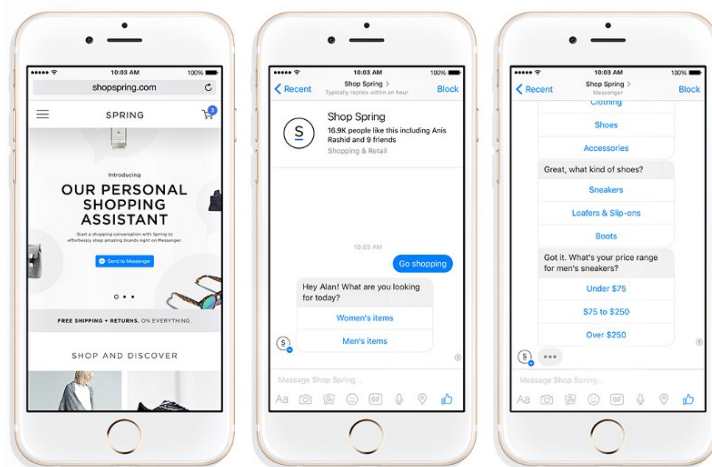


Figura 2.2: Cumpărături de haine folosind chatbot-ul companiei Spring pe Facebook Messenger

2.2 Modele generative

Spre deosebire de modelele anterioare, cele generative nu se bazează pe răspunsuri predefinite ci generează noi răspunsuri pornind de la zero. Modelele generative folosesc de obicei tehnici din Machine Translation, dar în loc de a traduce dintr-o limbă într-alta, vom ”traduce” de la o intrare la o ieșire (răspuns). Acestea oferă o mai bună impresie de comunicare cu un om real. Totuși, ele sunt extrem de greu de antrenat, sunt predispuse la erori gramaticale (în special unde lungimea propoziției este mai mare) și necesită o cantitate mare de date de antrenare.

Prezentul este încă sub semnul întrebării pentru acest tip de model, însă în următorii ani, acestea vor căpăta tot mai multă atenție și popularitate devenind tot mai performante. Dacă un chatbot va doborî testul Turing, șansele ca acesta să fie generativ sunt destul de mari. Deoarece modelele generative reprezintă o arie de cercetare încă nefinisată, acestea nu sunt folosite momentan în producție. În capitolul 3 vor fi prezentate diferite arhitecturi de rețele neurale, capabile să modeleze conversații.

Capitolul 3

Arhitectura

Modelarea unui limbaj reprezintă o sarcină extrem de dificilă pentru un calculator. Recente progrese în aria Deep Learning au făcut posibile diverse dezvoltări în această direcție, însă lucrurile sunt departe de a fi rezolvate. Arhitectura folosită pentru construcția chatbot-ului prezentat în această lucrare va fi expusă precum un bloc de construcții, plecând de la noțiunile de bază, până la arhitectura finală.

3.1 Rețele neurale artificiale

O rețea neurală artificială (RNA) reprezintă o paradigmă bazată pe procesare de informații a cărei inspirații provine din sistemul nervos biologic. Precum creierul uman, o RNA este compusă dintr-un număr mare de elemente de procesare interconectate (neuroni) lucrând la unison pentru a rezolva diverse probleme. RNA, precum oamenii, învață din exemple. Învățarea în sistemele biologice implică ajustarea conexiunilor sinaptice care există între neuroni. Acest principiu se aplică și acestor rețele.

Ca și modelare matematică propriu-zisă, o RNA poate fi observată în Figura 3.1. Avem o intrare reprezentată în figură de vectorul n -dimensional (x_1, x_2, x_3, x_4) . Această intrare reprezintă caracteristicile unui eșantion care face parte dintr-un set de date. Straturile intermediare se numesc straturi ascunse, iar rolul lor este să producă o abstractizare cât mai complexă a datelor de intrare, folosind diverse funcții cu activare neliniară. Ultimul strat se numește strat de ieșire, reprezentând eticheta (clasa) vectorului de intrare.

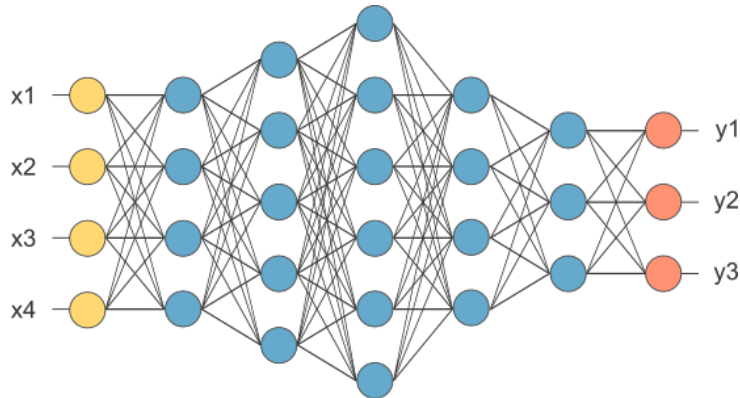


Figura 3.1: Exemplu de rețea neurală artificială

Ceea ce de fapt acest tip de rețele încearcă să învețe sunt ponderile dintre straturile sale. Se pleacă de la un set de ponderi alese aleator ¹ și se ajustează folosind algoritmul de propagare înapoi a erorii până când rețeaua se stabilizează.

TO CONTINUE

3.2 Rețele neurale artificiale recurente

După cum se poate observa, o RNA este utilă atunci când intrarea este formată dintr-un singur element, de exemplu o imagine. Modelarea unui limbaj însă presupune ca intrările să fie fraze. O frază este formată din mai multe elemente (cuvinte), deci o RNA nu poate modela o astfel de intrare. Soluția pentru această problemă este oferită de rețelele neurale artificiale recurente (RNAR). Acestea sunt folosite pentru a modela secvențe unde există o dependență temporală (fraze, serii de timp). Fiecare intrare curentă din secvență este condiționată de cele precedente. O astfel de rețea se poate observa în Figura 3.2. Straturile ascunse într-o RNAR au rolul de a păstra o captură a tot ceea ce s-a oferit ca intrare până în momentul curent. Valoarea fiecărui strat ascuns curent depinde astfel atât de intrarea curentă cât și de ieșirea stratului ascuns anterior. Putem considera toată această modelare ca pe o probabilitate condiționată $P(x_n|x_1, x_2, \dots, x_{n-1})$, unde în cazul unei fraze x_1, x_2, \dots, x_n reprezintă cuvintele acesteia.

¹Valorile ponderilor sunt de obicei subunitare iar inițializarea este făcută urmând diverse principii matematice bine definite. Pentru cazuri simpliste, inițializarea poate fi totuși făcută folosind o distribuție uniformă.

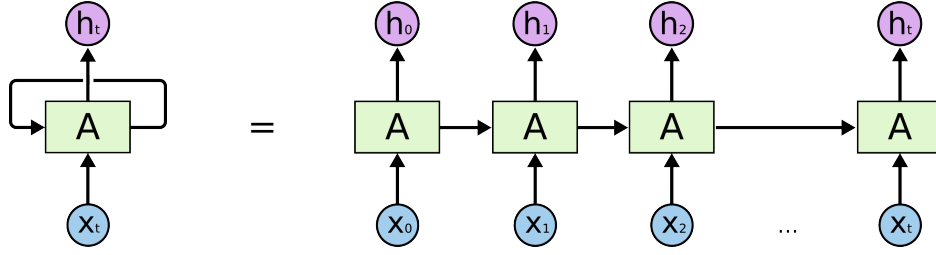


Figura 3.2: Exemplu de rețea neurală artificială recurentă

Algoritmul de reglare a ponderilor se numește propagarea înapoi în timp a erorii. O problemă serioasă care a apărut odată cu introducerea acestor rețele se numește risipirea gradientilor. În cazul în care secvența de intrare are o lungime mare, informația nu reușește să se propage în timp deoarece gradientul ponderilor se risipește, ponderile devenind 0. Mecanismul care înlătură această problemă se numește Long Short Term Memory (LSTM). Principiul este ca prin folosirea unor porți de transmitere a informației, gradientul să fie stabilizat. Acest mecanism a reprezentat un avans spectaculos pentru Deep Learning, permițând modelarea secvențială cu reținere a informației pe o perioadă îndelungată de timp. Ecuatiile 3.1 descriu calculele pentru porțile LSTM.

$$\begin{aligned}
 f_t &= \sigma(W_f x_t + U_f h_{t-1} + b_f) \\
 i_t &= \sigma(W_i x_t + U_i h_{t-1} + b_i) \\
 o_t &= \sigma(W_o x_t + U_o h_{t-1} + b_o) \\
 c_t &= f_t \circ c_{t-1} + i_t \circ \tanh(W_c x_t + U_c h_{t-1} + b_c) \\
 h_t &= o_t \circ \tanh(c_t)
 \end{aligned} \tag{3.1}$$

Poarta f_t se numește poartă de uitare (forget gate). Rolul ei este să stabilească câtă informație se va uita din trecut. Poarta i_t este poarta de intrare (input gate). Aceasta delimitează care este cantitatea de informație care se păstrează din intrarea curentă. Poarta o_t este cea de ieșire (output gate). Ea controlează ce informație va fi transmisă către ieșire. c_t se numește starea internă a celulei LSTM. Aceasta este calculată ca o combinație între starea anterioară a celulei, poarta de intrare și cea de ieșire. h_t reprezintă ieșirea curentă a rețelei și depinde de poarta de ieșire și starea celulei LSTM.

O RNAN aduce mai aproape ideea de modelare lingvistică, însă se poate observa că aceasta nu permite ca intrare decât o frază pe rând. Modelarea dorită este o pereche de forma întrebare-răspuns.

3.3 Modele Sequence-to-Sequence

Translația a reprezentat mereu un punct de interes în mediul procesării naturale de limbaj, constituind de altfel o provocare uriașă de-a lungul ultimilor ani. Până în anul 2014, majoritatea modelelor de translație se bazau pe lanțuri Markov ascunse (Hidden Markov Models - HMM), totul urmând a se schimba odată cu introducerea unei noi arhitecturi în acel an de către Cho et Al. Noul tip de arhitectură se numea Sequence-to-Sequence (seq2seq) și urma să aducă îmbunătățiri spectaculoase atât pe partea de translație cât și pe partea conversațională. Modelul este vizibil în Figura 3.3.

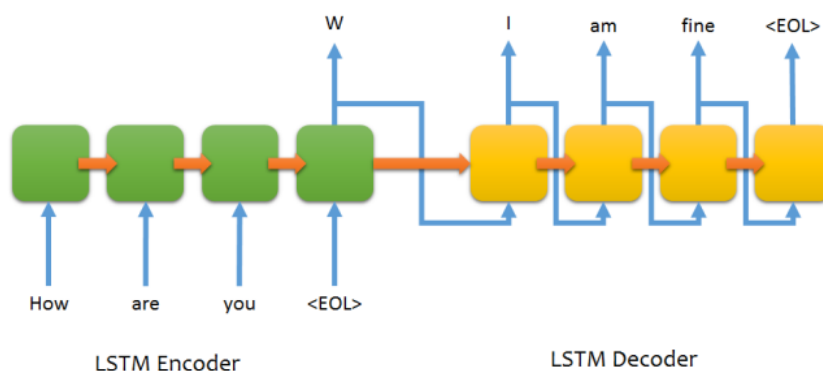


Figura 3.3: Modelul Sequence-to-Sequence

Aceast tip de arhitectură este împărțită în două jumătăți: codor (encoder) și decodor (decoder). Sarcina codorului este să primească o frază cu un număr variabil de cuvinte ca intrare iar unica sa ieșire ² să fie o captură a întregii intrări (un vector de lungime fixă). Această captură poate fi privită ca o sumarizare a întregii fraze. Sarcina decodurului este de a învăța fraze pornind de la ieșirea codorului, care va deveni prima intrare din decodor, și restul intrarilor precedente. Intrarea curentă w_n în decodor este ieșirea de la timpul anterior, w_{n-1} . Modelarea se transformă astfel într-o probabilitate condiționată: $P(w_n|w_1, w_2, \dots, w_{n-1}, c)$, unde c reprezintă ieșirea codorului, deseori întâlnit sub numele de context în literatura de specialitate.

Intrările în acest model sunt de forma întrebare-răspuns. Original folosit ca model de translație, unde intrarea pentru codor era fraza într-o limbă iar intrarea în decodor era fraza tradusă în limba dorită, acest principiu se

²Ieșirea ultimului element din secvență

poate aplica la fel de ușor pentru a modela o conversație. Spre deosebire de traducere unde totul se întâmplă punctual iar răspunsul nu depinde decât de intrarea curentă, o conversație este foarte dependentă de un context. Fără acest context, partenerul angrenat în discuție ar răspunde mereu luând în considerare doar ce i s-a spus în momentul de față, ignorând orice replică anterioară. Acesta nu este un comportament dorit în cadrul unei conversații și de aceea modelul seq2seq nu este destul de puternic de sine stătător pentru a modela o discuție.

3.4 Hierarchical Recurrent Encoder-Decoder

Modelarea contextului discuției reprezintă una dintre principalele nevoi în ceea ce privește o conversație care dorește să pară cât mai reală. Până recent, cea mai apropiată arhitectură care realiza acest lucru era modelul seq2seq, însă contextul era reținut doar la nivelul unei singure fraze. În anul 2016, Iulian Șerban et Al. a introdus rețeaua Hierarchical Recurrent Encoder-Decoder (HRED - Figura 3.4). Ea poate fi văzută ca o extensie peste seq2seq. Avantajul acesteia este că poate reține contextul discuției.

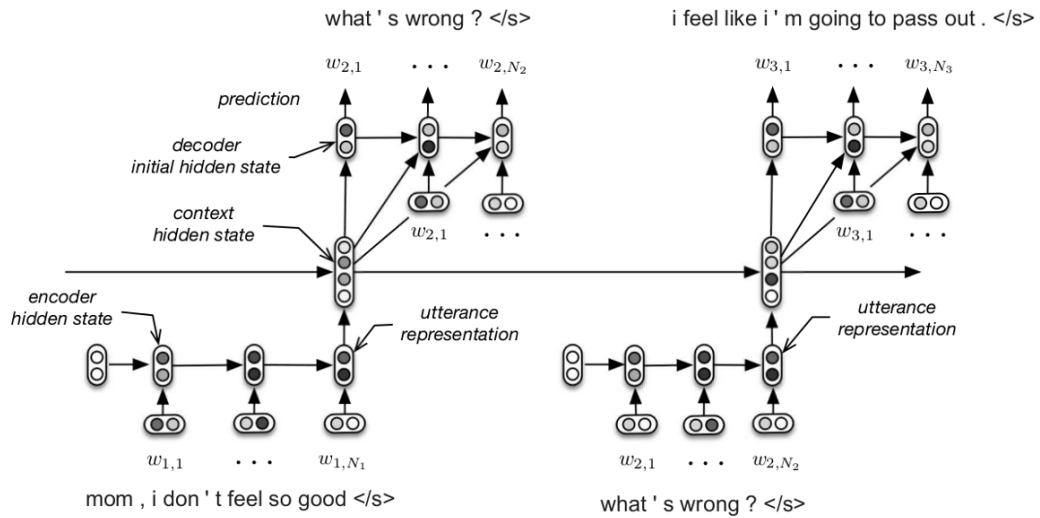


Figura 3.4: Hierarchical Recurrent Encoder-Decoder

După cum se poate observa, modelul este construit prin alăturarea mai multor rețele seq2seq legate între ele printr-o RNAR contextuală. După cum sugerează și numele, rolul acestei RNAR este de a reține contextul discuției de-a lungul mai multor fraze. Spre deosebire de seq2seq, ieșirea codorului nu

mai este oferită direct ca și primă intrare pentru decodor, aceasta trecând mai întâi prin RNAR contextuală. Pentru a fi mai ușor de înțeles ce se întâmplă, ne putem imagina toate cele n codoare și decodoare ca fiind simple unități de intrare, respectiv ieșire într-o RNAR obișnuită. Precum sunt într-o RNAR ieșirile dependente atât de intrarea curentă, cât și de cele precedente, prin analogie putem deduce că într-o rețea HRED, fiecare ieșire depinde de fraza curentă și de contextul discuției (frazele precedente).

3.4.1 Celula GRU

Gated Recurrent Unit sau GRU este un tip de celulă pentru o RNAR, menită să o înlocuiască pe cea LSTM. Ecuațiile din spatele celulei(3.2) sunt asemănătoare cu cele LSTM.

$$\begin{aligned} z &= \sigma(x_t U^z + s_{t-1} W^z) \\ r &= \sigma(x_t U^r + s_{t-1} W^r) \\ h &= \tanh(x_t U^h + (s_{t-1} \circ r) W^h) \\ s_t &= (1 - z) \circ h + z \circ s_{t-1} \end{aligned} \tag{3.2}$$

GRU are doar două porți: poarta r de resetare (reset gate) și poarta z de actualizare (update gate). Intuitiv, poarta de resetare determină cum se va combina noua intrare cu memoria precedentă iar cea de actualizare definește cât de multă memorie anterioară se păstrează. Principalele diferențe între GRU și LSTM sunt:

- GRU are doar două porți, pe când LSTM trei.
- GRU nu posedă o memorie internă c_t care să fie diferită de starea ascunsă expusă. Nu conține poarta de ieșire care este prezentă pentru LSTM.
- Poarta de intrare și cea de ieșire sunt cuplate de poarta de actualizare, iar poarta de resetare este aplicată direct stării ascunse anterioare. Astfel, responsabilitatea porții de resetare din LSTM este împărțită între cea de resetare și cea de actualizare de la GRU.
- Nu se aplică o a doua neliniaritate atunci când se calculează ieșirea.

Nu există o arhitectură câștigătoare clară între LSTM și GRU. Având mai puțini parametri (U și W), GRU se antrenează mai rapid și are nevoie de mai puține date pentru a generaliza. Pe de altă parte, o cantitate mare de date exprimă o putere mai mare de modelare din partea LSTM și astfel

se pot obține rezultate mai bune. Pentru realizarea acestei aplicații a fost aleasă o arhitectură de tip GRU.

3.4.2 HRED bidirecțional

Rolul codorului, precum a fost menționat mai devreme, este să captureze informația unei fraze într-un vector de lungime fixă. În orice limbaj, sensul unui cuvânt nu este stabilit doar de cuvintele precedente ci și de cele ce vor urma. Deoarece RNAR modelează cuvintele dintr-o frază pornind de la cuvântul w_1 la w_n , înțelesul din viitor al acestora este ignorat, iar captura frazei poate să nu fie îndeajuns de reprezentativă. Astfel, este propusă folosirea unei RNAR bidirecționale (Figura 3.5) pentru codor. Acest tip de rețea folosește două parcurgeri ale secvenței de intrare: cea înainte, care se desfășoară în mod obișnuit și cea inversă (de la w_n la w_1), pentru capturarea înțelesului din viitor al cuvintelor. Ieșirea pentru parcurgerea înainte va fi la timpul n iar cea pentru parcurgerea inversă va fi la timpul inițial. Vectorul de lungime fixă va fi format din concatenarea celor două ieșiri ale rețelei bidirecționale.

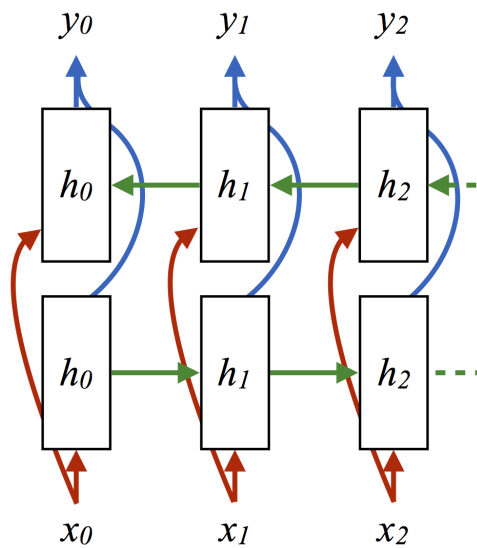


Figura 3.5: RNAR bidirecțional

3.5 Word Embeddings

TODO

3.6 Predicție

Având modelul antrenat, acesta trebuie folosit pentru predicție, adică generarea unor răspunsuri adecvate pentru contextul discuției. Generarea cuvintelor este condiționată atât de context, cât și de cuvintele anterior generate. Spre deosebire de antrenare, în momentul predicției, secvența de ieșire este goală inițial (un vector de zerouri). Cuvintele sunt generate secvențial, unul câte unul și adăugate pe poziția t la care a ajuns secvența. Deoarece stratul softmax al secvenței de ieșire returnează un vector de probabilități peste toate cuvintele din vocabular, pentru a returna cuvântul dorit se alege cel cu probabilitatea cea mai mare. Pornind de la această idee, există două variante de generare: metoda greedy și beam search.

3.6.1 Greedy

Precum în teoria clasică a algoritmicii, scopul metodelor greedy este să selecteze întotdeauna optimul local. În cazul de față, optimul local reprezintă cuvântul cu probabilitatea de apariție cea mai mare. Prin această metodă nu este garantat la final că fraza produsă este cea mai bună cu putință, deoarece o altă combinație de cuvinte poate produce oricând o frază cu un scor mai bun³. Astfel, singurul avantaj al acestei metode este viteza de predicție rapidă. În practică, se preferă folosirea unor algoritmi mai complexi, capabili să aleagă fraza cea mai potrivită dintr-un set de fraze candidat.

3.6.2 Beam search

Algoritmul beam search folosește ideea de arbore de căutare pentru a genera fraze (Figura 3.6). În locul alegerii cuvântului cu probabilitatea cea mai mare la fiecare pas, se vor alege top k cuvinte cu cele mai mari probabilități. În literatură, k se mai numește și beam size. Fiecare nod (cuvânt) ales va produce un număr n_c de copii. În felul acesta este garantată generarea mai multor fraze candidat din care se va alege. Fiecare frază are un scor $S(f)$, care este de forma ecuației 3.3, unde $P(w_i)$ reprezintă probabilitatea cuvântului ales. O frază devine candidat atunci când cuvântul curent generat este simbolul de sfârșit de propoziție. În momentul în care algoritmul a terminat de generat toate frazele candidat, cea mai bună este determinată ca fiind cea cu scorul $S(f)$ cel mai mare.

³Sumă din logaritmul probabilităților frazei

$$S(f) = \frac{1}{n} \sum_{i=0}^n \log P(w_i) \quad (3.3)$$

Pentru ca acest arbore de căutare să nu capete o creștere exponențială, la fiecare pas se păstrează un număr egal cu beam size cele mai bune fraze iar restul sunt decartate. Acest algoritm nu asigură de departe optimul global, însă spațiul căutării oferit de folosirea unui arbore este mult mai bine dezvoltat și ales decât la metoda greedy. Cu cât beam size și numărul de copii generați sunt mai mari, cu atât algoritmul oferă soluții cât mai apropiate de optimul global. De asemenea, odată cu creșterea acestor parametri, crește și timpul de căutare și generare în arbore. De aceea, trebuie să existe un compromis între calitatea predicției și viteza de execuție a algoritmului.

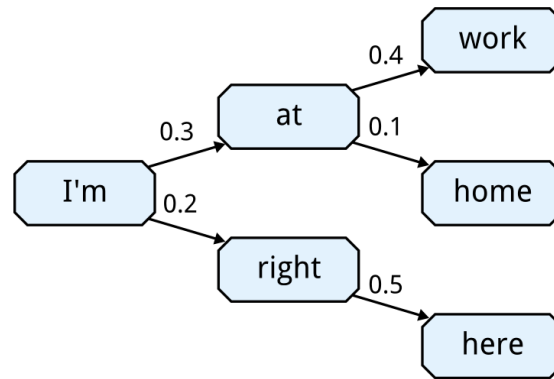


Figura 3.6: Arbore de căutare pentru algoritmul beam search