

UNIVERSITATEA "ALEXANDRU IOAN CUZA" DIN IAI
Facultatea de Informatica din Iasi

LUCRARE DE LICENTA

CREARE RETELELOR DE SORTARE FOLOSIND
DEEPLARNING

DANILA MARIUS CRISTIAN

Data: *Iulie, 2018*

Coordonator
Lect. Dr. Frasinaru Cristian

Abstract

Scopul acestei lucrari este de a construi un algoritm de deeplearning care sa genereze compartorii dintr-o retea de marime N . N reprezinta numarul de fire din retea. Desi exista abordari neimperative functionale, care genereaza rezultate intr-un timp multumitor, nu s-a mai incercat pana acum generarea de retele de sortare prin aceasta abordare. Rezultatul obtinut este un model antrenat care produce rezultate pentru retele de dimensiuni mici.

Cuprins

1	Introducere	3
2	Descrierea solutiei	4
2.1	Arhitectura proiectului	4
2.2	Arhitectura modelului	7
2.3	Antrenarea modelului	8
3	Retele de sortare	9
3.1	Prezentare	9
3.2	Generarea retelelor de sortare folosind algoritmi evolutivi	11
3.3	Retelele optimale	11
3.3.1	Prezentare	11
3.4	Principiul 0-1	12
3.4.1	Enunt	12
3.4.2	Demonstratie	12
3.5	Metode imperative de generare a retelelor de sortare .	12
3.6	Aplicatii	14
4	Retele neuronale, prezentare	14
4.1	Arhitecturi de retele neuronale	14
4.2	Neural Turing Machine (NTM)	15
4.2.1	Descriere	15
4.2.2	Aplicatii	17
5	Tehnologii folosite	17
5.1	NTMLasagne	17
5.2	Lasagne	18
5.3	Theano	18
5.4	Python	18

5.5	Java	18
5.6	Spring Boot	19
5.7	RaphaelJS	19
5.8	Twitter Bootstrap	19
6	Concluzii	20

1 Introducere

O retea de sortare reprezinta un model abstract din informatica, compus din fire si comparatori. Scopul comparatorilor este de a interschimba 2 numere de pe firele pe care sunt plasati. Desi sunt modele simple, retelele de sortare au aplicatii numeroase in calculul paralel si sunt folosite la sortarea seturilor de numere in cadrul placilor grafice.

In cadrul acestei lucrari mi-am propus sa dezvolt un model capabil sa genereze, pe baza unui numar N de fire primit ca input, lista de comparatori necesari pentru o retea de sortare cu N fire. In cadrul proiectului am dezvoltat si o interfata web ce permite vizualizarea cat si simularea unei sortari pe reteaua obtinuta de model. Cu aceasta ocazie scoatem in lumina si elementul de inovatie al acestei lucrari, si anume ca pana acum nu s-a mai incercat generarea de retele de sortare prin aceasta abordare. In spate, modelul foloseste o arhitectura de retele neuronale introdusa recent (*20 Oct 2014*), optimizata pentru invatarea de algoritmi imperativi. Modelul obtinut produce rezultate pentru retele de dimensiuni mici, de 2, 3 sau 4 fire. Nu produce comparatori pentru retele care sa sorteze complet siruri de numere mai mari de 5 elemente, dar acest lucru lasa loc de extindere a proiectului pe viitor.

Motivatia din spatele alegerii proiectului de fata este data de viteza de generare ca timp a retelelor de sortare ce este semnificativ mai buna pe retele de dimensiuni mari fata de celelalte abordari.

Proiectul inglobeaza 2 modele antrenate pe seturi de antrenament diferite, bazate pe rezultatele a 2 algoritmi imperativi. Acesti 2 algoritmi imperativi au la baza modul de generare a interschimbarilor din algoritmi de sortare *Bubble Sort* si *Insertion Sort*. Fiecare model incearca sa produca rezultatele similare cu unul dintre cei 2 algoritmi imperativi folositi. Cele doua modele sunt rulate de o a doua componenta a proiectului, *model manager*, intr-un thread. Scopul acestei componente este sa centralizeze rezultatele de la cele doua modele si sa comunice cu componenta *backend* din spatele interfetei web. Intr-un final, ultima componenta a proiectului este cea de interfata web care afiseaza retelele produse de cele doua modele.

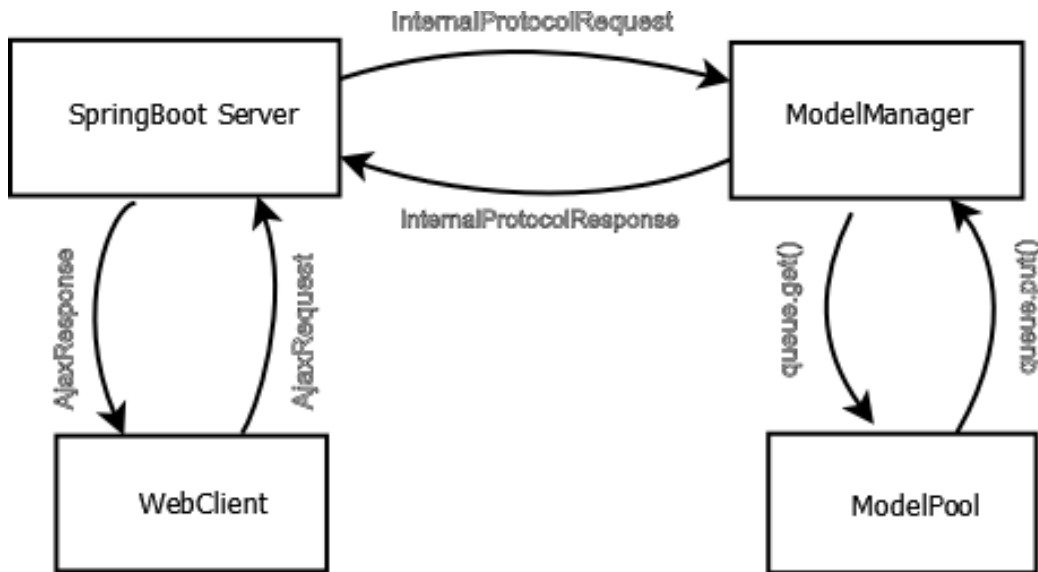
2 Descrierea solutiei

Solutia problemei consta intr-un model antrenat pe un set de date produs de un algoritm imperativ de generare a retelelor de tip bubble sorting network si insertion sorting network. Peste modele generate avem o interfata web se permite vizualizarea in browser a retelei generate.

2.1 Arhitectura proiectului

Proiectul complet este compus din 4 componente:

- Un server web, bazat pe SpringBoot. Rolul serverului web este sa-i livreze utilizatorului codul interfetei web a aplicatiei. De asemenea mediaza requesturile dintre requesturile unui utilizator de a primi comparatorii pentru un anumit numar de fire si *model manager*.
- O interfata in browser. Rolul interfetei in browser este acela de a se putea vizualiza efectiv reseaua generata de un modelul selectat pe un anumit numar de fire. Utilizatorul alege numarul de fire si modelul pe baza caruia se genereaza outputul.
- O componenta numita model manager. Managerul este compus dintr-un server TCP/IP, prin care care primeste requesturi de la serverul web. O clasa care proceseaza requestul primit de la serverul web si in functie de parametrii din request cere output de la thread pool-ul care ruleaza modelele.
- Un thread pool, fiecare model antrenat ruleaza intr-un thread separat din acest pool. Aceasta componenta incarca modelul antrenat de pe disc si ruleaza intr-un thread separat functia de predictie. Motivatia acestei abordari este legata de directiile de extindere pe viitor a proiectului.



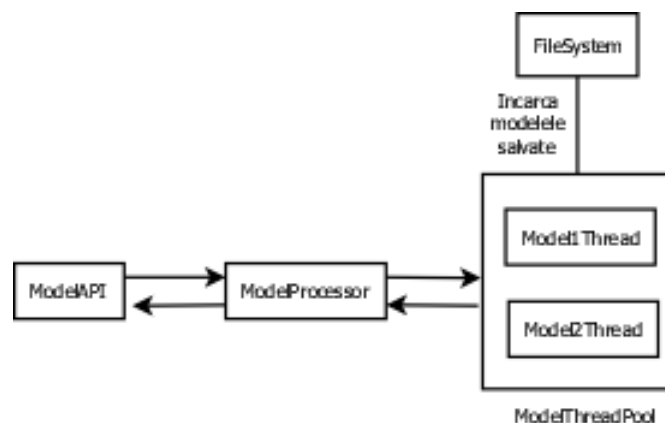
Arhitectura generara a proiectului

In diagrama de mai sus observam schematic, abstractizat, principalele componente ale proiectului si relatiile dintre ele. In continuare descriem in detaliu aceste relatii:

- **SpringBoot Server - WebUI**, etapa principala din ciclul de viata al proiectului este livrarea codului interfetei clientului, apoi, de interes este interactiunea dintre *WebUI* si server. De obicei fluxul normal este:
 1. Utilizatorul selecteaza numarul de fire, modelul antrenat folosit si asteapta producerea rezultatului. In spate, *WebUI* impacheteaza optiunile utilizatorului intr-un pachet *JSON* si trimite prin *AJAX* un request catre server.
 2. *Serverul* primeste requestul de la *WebUI* printr-un endpoint. In functie de optiunile alese, se creeaza un apel catre clasa *ManagerClient* care are rolul de a fi un adaptor intre server si *ModelManager*.
 3. *Server-ul* impacheteaza intr-un *JSON* lista de comparatori returnata de *ManagerClient* si il trimite *WebUI*-ului.
 4. *WebUI-ul* primeste *response*-ul *AJAX* cu lista de comparatori si initiaza un apel spre modulul de desenare, care creeaza reprezentarea grafica a retelei produse.

- **SpringBoot Server - ModelManager**, comunicarea dintre aceste componente este, de regula, initializata de *ManagerClient* si cronologic este structurata astfel:
 1. *ManagerClient*-ul initiaza o conexiune TCP/IP cu *ModelManager*-ul. Ii trimite acestuia un pachet codificat intr-un protocol propriu ce contine numarul de fire si modelul pe baza caruia sa se genereze comparatorii.
 2. *ModelManager*-ul accepta conexiunea de la *ManagerClient*, primeste pachetul, il decodifica si realizeaza un apel catre *ModelThreadPool*, apoi impacheteaza lista de comparatorii generata si o trimite inapoi la *ManagerClient*.
- **ModelManager - ModelPool**. De la initializare *ModelPool*-ul incarca de pe disc fiecare dintre modele si porneste un thread separat ce le ruleaza.
 1. In functie de modelul ales *ModelManager*-ul realizeaza un apel la metoda respectiva din *ModelPool*.
 2. Prin intermediul unei cozi *thread-safe*, *ModelPool*, semnalizeaza threadul cu modelul respectiv ca trebuie evaluat un nou rezultat. Rezultatul modelului este transmis inapoi prin acelasi mecanism.

In continuare dorim sa analizam mai atent structura si modul de functionare al *ModelManager*-ului.



Arhitectura interna a model managerului

În diagrama de mai sus sunt expuse principalele submodule ale managerului. Expunem succint scopul fiecărui submodule:

- *ModelApi*, reprezintă un server TCP/IP ce rulează pe portul 2018. Comunicarea cu componenta de *ModelManager* se face prin trimiterea de TCP/IP pe acest port.
- *ModelProcessor*, decodifică pachetele primite de eventualii utilizatori apelează metoda potrivită din urmaorul submodule. *ModelProcessor*-ul are și rolul de a aproxima probabilitățile din lista de comparatori generată de model.
- *ModelThreadPool*, încarcă de pe disc fiecare dintre modelele salvate și apoi creează un thread în care rulează fiecare model. Comunicarea dintre threadul principal se face prin intermediul a două cozi *thread-safe*. Threadul principal pune în coada numărul de fire dorit, iar threadul ce rulează modelul respectiv pune în a doua coadă lista de comparatori generată.

2.2 Arhitectura modelului

Modelul este un "Neural Turing Machine". Pentru construcția modelului am folosit în spate un framework numit ntm-lasage, optimizat pentru crearea de modele de acest fel. Se abstractizează în acest fel detalii de implementare.

Modelul conține o memorie cu un layout de 1024x160 celule. Ca funcție de loss am folosit crossentropia binară, pentru că intern numerele ce vin ca input în model sunt reprezentate ca siruri binare.

Procesul de backpropagation este abstractizat de model. Actualizarea weighturilor din rețea este optimizată cu rmsprop. Datele de antrenament sunt procesate de rețea în batchuri de 1.

Structura inputului rețelei este de un vector de 9 elemente ce reprezintă codificarea în binar a numărului de fire din rețea. Outputul rețelei este un vector tridimensional cu $2 * \text{numar_comparatori}$ elemente. Elementele din vector sunt semantic grupate 2 câte 2, codificate ca un sir de biti, și reprezintă pozițiile comparatorilor din rețea.

2.3 Antrenarea modelului

Modelul ce genereaza retele bubble sorter a fost antrenat pe un set de input 1000 de elemente, pe numere mai mari de atat timpul necesar finalizarii procesului de antrenament devine foarte mare. Datele din setul de antrenament au fost generate conform algoritmului:

```
def bubble_sorter(n):
    if n <= 1:
        return []

    ret_val = []
    for i in range(0, n - 1):
        ret_val.append((i, i + 1))
    ret_val += self.dummy_sorter(n - 1)
    return ret_val
```

si arata astfel pentru un set de antrenament de 5 elemente:

```
[2, [(0, 1)]]
[3, [(0, 1), (1, 2), (0, 1)]]
[4, [(0, 1), (1, 2), (2, 3), (0, 1),
      (1, 2), (0, 1)]]
[5, [(0, 1), (1, 2), (2, 3), (3, 4),
      (0, 1), (1, 2), (2, 3), (0, 1), (1, 2), (0, 1)]]
[6, [(0, 1), (1, 2), (2, 3), (3, 4),
      (4, 5), (0, 1), (1, 2), (2, 3),
      (3, 4), (0, 1), (1, 2), (2, 3), (0, 1), (1, 2), (0, 1)]]
```

Insa date de mai sus reprezinta la nivel abstract inputul si outputul dorit pentru retea. In proiect datele sunt reprezentate in binar ca un sir de 9 biti, iar interschimbarile nu sunt reprezentate explicit ca tuple. Reteaua va produce, practic, o lista de siruri biti ce vor fi grupate, in ordine, 2 cate 2.

Spre exemplu, pentru urmatorul input:

```
numpy.array([[0, 0, 0, 0, 0, 0, 0, 1, 1]])
```

reprezinta dimensiunea 3, adica retea va trebui sa genereze setul de comparatori pentru o retea de 3 elemente. Inputul de mai sus a fost testat pe o retea antrenata pe un set de date 1000 elemente, modelul este Bubble Sorter. Outputul produs este urmatorul:

```
[[[0.000751546365971, 0.465849417004],
  0.000769892576671, 9.00920625725e-05,
  0.000637261131599, 7.82020981282e-05,
  0.000600783421602, 6.5095684253e-05,
  0.12773718494, 6.73165195701e-05,
  0.239040904467, 0.180481531645,
  0.214272843324, 0.244833481652,
  0.370771723772,
```

0.273710496834 ,	0.510196603956] ,
0.454686659751 ,	
0.496889833317] ,	
[4.46837405637e-05 ,	[3.10455175941e-05 ,
3.44817119736e-05 ,	2.33229850784e-05 ,
2.6885035585e-05 ,	1.51453999616e-05 ,
3.21667374046e-05 ,	2.18735040572e-05 ,
0.262205621903 ,	0.376419535208 ,
0.240104608732 ,	0.237146807326 ,
0.341346314965 ,	0.429571992622 ,
0.501759117299 ,	0.537611938861 ,
0.507289515972] ,	0.511490463268] ,
	[2.89867797738e-05 ,
[3.48919682303e-05 ,	2.23753419482e-05 ,
2.59198402876e-05 ,	1.35824510377e-05 ,
1.83589819237e-05 ,	2.04069513415e-05 ,
2.46684454266e-05 ,	0.406359755172 ,
0.330040528556 ,	0.238643669264 ,
0.237285575094 ,	0.452255805808 ,
0.394284580217 ,	0.54490706061 ,
0.525120693702 ,	0.513243797725]]]

Semnificatia fiecarui double din lista este probabilitatea ca bitul respectiv sa fie 1. Acuratetea probabilitatilor se sporeste cu antrenarea sporita a modelului. Intr-un final, se grupeaza 2 cate 2 elementele din lista si se decodifica din binar in decimal.

2.4 Prezentarea interfetei

– pune 1-2-3 screenshoturi – vorbește foarte succint despre fiecare optiune pe care o are utilizatorul

3 Retele de sortare

In cadrul acestui capitol prezentam acestui capitol prezentam fundamentele teoretice necesare pentru intelegerea conceptului de *retea de sortare* cat si aplicatiile acestora in practica

3.1 Prezentare

O retea de sortare cu n inputuri este o secventa fixa de comparatii si inter-schimbari (comparatori) ce sorteaza toate inputurile de marime n . Deoarece

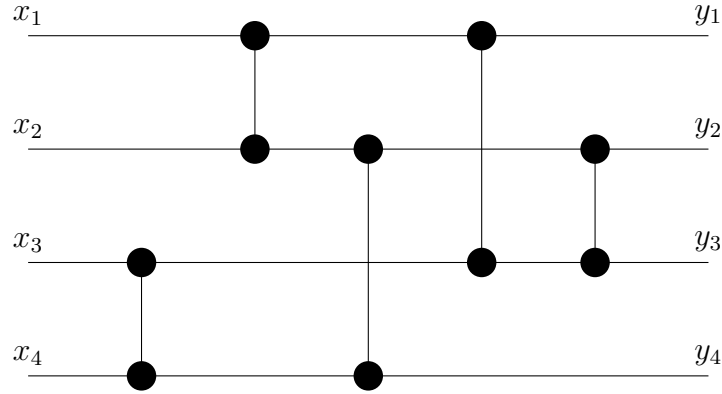


Fig. 1: O retea de sortare cu 4 inputuri. Valorile de input x_1, x_2, x_3, x_4 din partea stanga a liniilor orizontale trec printr-o secventa de operatii de comparatii si interschimbari, reprezentate prin linii verticale ce conecteaza perechi de linii orizontale. Fiecare comparator de acest fel isi sorteaza cele doua valori de input, rezultand intr-un final in liniile orizontale ce contin valorile sortate de output $y_1 \leq y_2 \leq y_3 \leq y_4$. Aceasta este o retea minimala din punct de vedere a numarului de comparatori folositi.

aceeasi secventa de comparatori sorteaza orice secventa de marime n , reprezinta un algoritm de sortare independent de date. Asta inseamna ca secventa de comparatii ce se executa nu depinde de inputul primit. Datorita structurii pe care o are o retea de sortare, sunt de preferat in implementarile paralele ale algoritmilor de sortare, precum cele de pe placile grafice.

Motivati de aceste aplicatii in practice, retelele de sortare sunt un subiect important de cercetare inca din 1950. Un interes deosebit il au retelele de sortare optimale, care folosesc numarul minim posibil de comparatori. Generarea retelelor de sortare optimale reprezinta o problema grea de optimizare, investigata pentru prima data de O'Connor si Nelson pentru $4 \leq n \leq 8$. Retelele gasite de cei 2 aveau numarul minim de comparatori pentru 4, 5, 6 si 8 inputuri, dar aveau nevoie de 2 comparatori in plus pentru 7 inputuri. Acest rezultat a fost imbunatatit in 1968 de Batcher care a gasit retele minimale pentru $n \leq 8$.

n	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
C	0	1	3	5	9	12	16	19	25	29	35	39	45	51	56	60

Fig. 2: Numarul minim de comparatori cunoscut pentru retele de sortare de dimensiune $n \leq 16$. Aceste retele au fost studiate intens, dar s-a dovedit ca doar rezultatele pentru $n \leq 8$ sunt minimale.

Pana in 2013, se cunosteau retele de sortare minimale doar pentru inputuri $n \leq 8$ (vezi Tabelul 1). In 2013, Valsalam a introdus un articol in care prezinta modalitati de generare a retelelor optimale prin algoritmi evolutivi si reuseste sa obtina noi rezultate pentru n de marime 17 , 18, , 19 20, , 21 si 22.

3.2 Generarea retelelor de sortare folosind algoritmi evolutivi

– descrierea succinta a abordarii prin algoritmi evolutivi

3.3 Retelele optimale

3.3.1 Prezentare

In aceasta sectiune a lucrarii dorim sa amintim despre *retelele lui Green*, intrucat folosind metoda lui *Green* de construire a retelelor obtinem numarul minimal de comparatori.

Mai jos afisam rezultatul metodei lui *Green* pentru o retea de marime 16:

Acesta retea are 60 de comparatori, care reprezinta cea mai mica valoare cunoscuta pentru o retea cu 16 inputuri[16][17]. Comparatorii dintr-o astfel de retea sunt simetric aranjati de la axa orizontala pana in mijlocul retelei. Acest tip de retele au reprezentat pentru unii cercetatorii obiectiv de studiu pentru dezvoltarea algoritmilor evolutivi care genereaza retele de sortare.

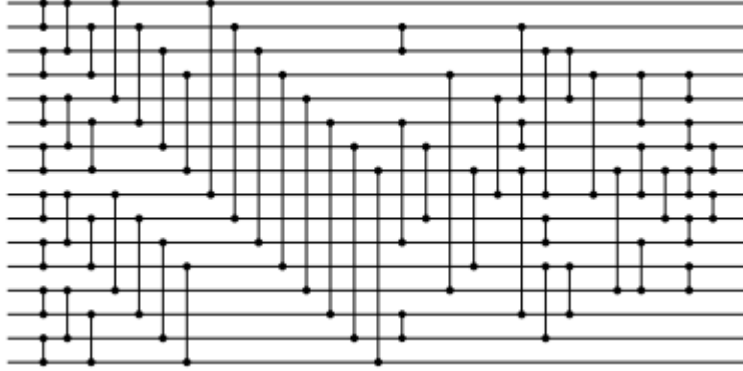


Fig. 3: Reteaua de comparatori a lui green pentru un input $n = 16$

3.4 Principiul 0-1

3.4.1 Enunt

Principiul 0-1 este o modalitate de verificare a corectitudinii unei rețele de sortare. Conform acestuia, o rețea de sortare care sortează corect orice secvență de 0 și 1, va sorta corect orice secvență de numere.

3.4.2 Demonstratie

Lemma: Presupunem f o funcție monotona, crescătoare. Atunci, dacă o rețea mapează x_1, x_2, \dots, x_n la y_1, y_2, \dots, y_n , va mapa $f(x_1), f(x_2), \dots, f(x_n)$ la $f(y_1), f(y_2), \dots, f(y_n)$.

Demonstratie: Prin inducție la numărul de comparatori din rețea folosind:

$$\begin{aligned} f(\min(a, b)) &= \min(f(a), f(b)) \\ f(\max(a, b)) &= \max(f(a), f(b)) \end{aligned}$$

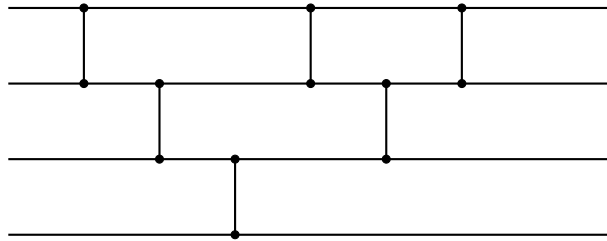
Să presupunem că o rețea nu este rețea de sortare. Atunci, va mapa valorile x_1, x_2, \dots, x_n arbitrare, la y_1, y_2, \dots, y_n . Unde $y_i > y_{i+1}$ pentru $1 \leq i < n$. Presupunem $f(x) = 1$ dacă și numai dacă $x \geq y_i$, 0 altfel. Rețeaua va mapa $f(x_1), f(x_2), \dots, f(x_n)$ la $f(y_1), \dots, f(y_i) = 1, f(y_{i+1}) = 0, \dots, f(y_n)$. Astfel, rețeaua nu va sorta toate inputurile de tipul 0 – 1.

3.5 Metode imperative de generare a retelelor de sortare

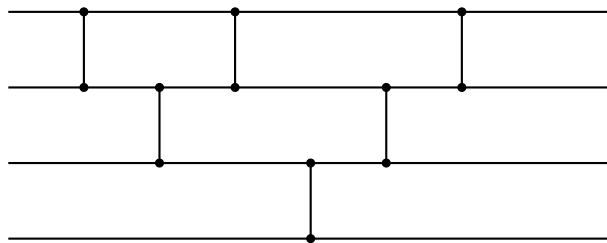
Exista mai multe metode imperative pentru a genera o retea de sortare. Cele mai usoare modalitati sunt de a genera o retea bazata pe algoritmi de sortare bubble sort si insertion sort. In cazul primei optiuni se adauga un comparator intre firele 0 si 1, 1 si 2, ..., $n - 1$ si n . Se repeta procesul de creare a comparatorilor inserand un comparator intre 0 si 1, ..., $n - 2$ si $n - 1$. Se repeta aceasta procedura de creare a comparatorilor pana la pasul in care inseram un singur comparator intre firele 0 si 1.

In cazul unei insertion sorting network, o retea bazata pe insertion sort procesul este similar, singura diferenta este ca buclele de creare a comparatorilor sunt in ordine inversa fata de metoda precedenta. Astfel, prima data se creeaza un comparator intre firul 0 si 1. Se creeaza un comparator intre 0 si 1, 2. Pana cand se ajunge la pasul in care se genereaza un comparator intre 0 si 1, ..., $n - 1$ si n .

Un exemplu de bubble sorting network:



Un exemplu de insertion sorting network:



Algoritmul imperativ pe care l-am folosit pentru a genera un bubble sorting network este:

```
def bubble_sorter(n):  
    if n <= 1:  
        return []
```

```

ret_val = []
for i in range(0, n - 1):
    ret_val.append((i, i + 1))
ret_val += self.dummy_sorter(n - 1)
return ret_val

```

Pentru un numar de fire dat, n , generam lista de comparatori, adica conexiunile dintre firele din retea.

Aceste arhitecturi de retele de sortare nu sunt optimale si vor sorta un sir de numere ca input in timp $O(n^2)$. Retelele optimale si cele mai des folosite in practica sunt bazate pe algoritmul bitonic sort, acestea ating complexitate de $O(n \log^2(n))$. Insa motivul pentru care am prezentat bubble sorting network si insertion sorting network este ca bazandu-ne pe acesti algoritmi am generat datele de antrenament pentru modelul dezvoltat. Motivatia ar fi ca, cel putin din experimentele realizate in cadrul acestui proiect, este mai usor de antrenat un model care sa produca date relevante atunci cand primeste date de antrenament generate pe baza acestor algoritmi.

3.6 Aplicatii

Cea mai la indemana aplicatie in practica a retelelor de sortare o reprezinta sortarea pe placile grafice. Placile grafice dispozitive orientate puternic spre paralelizarea instructiunilor, pe arhitectura *single-instruction, multiple data* (SIMD). Si avand in vedere si ca placile grafice depasesc in performanta CPU-ul cand vine vorba de algoritmi cu limita de memorie si timp de calcul, retelele de sortare devin o optiune eficienta.

4 Retele neuronale, prezentare

Retelele neuronale reprezinta o ramura a inteligentei artificiale. Sunt modele abstracte ce incearca sa simuleze creierul uman in modul de intercationare cu natura. O retea neuronală este compusa din neuroni artificiali. Neuronii artificiali din retea sunt dispusi pe mai multe straturi. Din exterior le putem vedea ca pe un black box conectat la un strat de neuroni de input, responsabil de preluarea inputului si la un strat de neuroni de output, responsabil pentru a genera outputul. In functie de neuronii din stratul de output care se vor activa se decodifica rezultatul produs de retea.

4.1 Arhitecturi de rețele neuronale

Retelele neuronale convolutionale sunt modele simple, proiectate să primească un set de input fix și să producă un set output de mărime fixă. În practică, acest tip de rețele sunt ideale pentru computer vision. Sunt mult mai eficiente din punct de vedere al resurselor consumate la sarcini precum recunoașterea de obiecte în imagini.

Retelele neuronale recurente se deosebesc de cele convolutionale prin posibilitatea de a primi un set de input variabil și a produce un set de output variabil. În practică se folosesc la procesarea limbajului natural. Acest tip de rețele neuronale sunt ideale și pentru rezolvarea problemei de generare a rețelilor de sortare, întrucât putem codifica numărul de fire din rețea ca o secvență de lungime variabilă, iar rezultatul, respectiv setul de comparatori pentru numărul de fire dat poate varia. Pe lângă aceste considerente, acest tip de rețele conține și o memorie asociată.

Un caz particular de rețele neuronale recurente sunt *masinile turing neurale*^{*}. Acestea fiind optimizate în a învăța algoritmi imperativi simpli. Un exemplu ar fi problema sortării, rezolvată în mod tradițional cu un algoritm imperativ poate fi rezolvată folosind acest tip de rețele. Pentru aceasta, rețeaua va primi ca input tuple de forma (X, Y) . Unde X reprezintă un sir de numere aflate în ordine random, iar Y reprezintă varianta sortată a sirului X cu un algoritm imperativ de sortare. Odată antrenată pe un set suficient de mare de date de acest fel, modelul va învăța practic să aplice algoritmul de sortare folosit pe componenta Y .

Motivul prezentării scenariului de mai sus este că hiperparametrii și structura straturilor de neuroni dintr-o rețea de neuronale folosită la sortarea de numere sunt similare în proporție foarte mare cu cei din rețeaua proiectată să genereze rețele de sortare.

4.2 Neural Turing Machine (NTM)

Retelele NTM fac parte din clasa rețelilor neuronale recurente și sunt introduse relativ recent. Au fost introduse într-un articol publicat de Alex Graves în 2014. Acest model de rețea neuronală combină capacitatea de pattern matching a rețelilor neuronale cu capacitatea algoritmică a calculatoarelor programabile. Mai exact, putem folosi acest model de rețea pentru a învăța algoritmi simpli. În paperul introdus de Graves sunt aduse rezultate pentru algoritmi de copiere a unui sir sau sortare. Aceste aspecte ne inspiră o

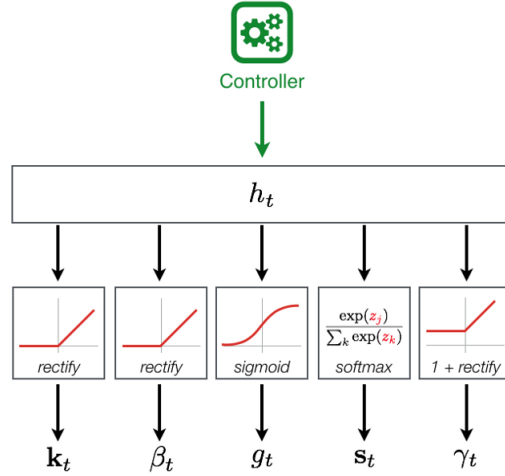


Fig. 5: Arhitectura capului de scriere/citire a unui *Neural Turing Machine*. Sursa: <https://medium.com/snips-ai/ntm-lasagne-a-library-for-neural-turing-machines-in-lasagne-2cdce6837315>

care este folosit de capetele de scriere si citire ca sa poata interactiona cu memoria. Reprezentarea aceasta nu este identica cu cea stocata in memorie. Tipul de controler pe care il alegem reprezinta cea mai semnificativa alegere arhitecturala. Controllerul poate fi fie o retea feed-forward, fie o retea recurenta.

Capetele de scriere si citire sunt singurele componente din cadrul unui Neural Turing Machine care interactioneaza in mod direct cu memoria. Intern, comportamentul fiecarui dintre capete este controlat de vectorul de weighturi care este actualizat la fiecare pas in timp. Fiecare weight din vector corespunde cu gradul de interactiune cu fiecare celula din memorie. Un weight de 1 focuseaza atentia NTM-ului spre celula respectiva.

Reprezentare vizuala a notiunilor descrise mai sus:

4.2.2 Aplicatii

– descriere aplicatii algoritmi paraleli de sortare – utilizare pe GPU – utilizare pe CPU

5 Tehnologii folosite

In cadrul proiectului am facut uz de o serie de frameworkuri si biblioteci pentru a fi scutiti de unele detalii de implementare. In aceasta sectiune vom prezenta pe scurt tehnologiile principale folosite.

5.1 NTMLasagne

NTMLasagne este frameworkul pe care l-am folosit in proiect ca sa creem modelul. Pe scurt, este un wrapper peste Lasagne ce ne permite sa creem "Neural Turing Machines" intr-un mod simplificat, abstractizand detalii de implementare. Ne sunt puse module pentru stratul NTM (en. NTMLayer) unde toate componentele precum capetele de scriere, controllerul si memoria pot fi customizate. In cadrul acestui proiect am ales sa nu customizam niciuna dintre componente. In schimb, le-am folosit pe cele default oferite de framework.

Motivatia din spatele alegerii acestui framework este, chiar daca documentatia este relativ putina, comunitatea puternica din spatele proiectului si faptul ca procesul de dezvoltare este activ.

5.2 Lasagne

Lasagne este o biblioteca minimalista folosita pentru construirea de retele neuronale. Permite construire de retele feed-forward, cat si de retele convolutionale, ofera implementari pentru metode de optimizare folosite des (i.e. ADAM si RMSProp). De asemenea, modele construite in Lasagne pot fi antrenate fie pe CPU sau GPU pentru ca biblioteca foloseste Theano in spate.

5.3 Theano

Theano este o biblioteca de python ce permite definirea si optimizarea de expresii matematice ce contin array-uri multidimensionale intr-un mod eficient. Chiar daca in momentul de fata Tensorflow este alegerea facuta de majoritatea proiectelor, theano este folosit de cele 2 biblioteci specificate anterior si este in general mai flexibil.

5.4 Python

Python este un limbaj de scripting, multi paradigma. Motivul pentru care am ales sa implementam modelul in python este sintaxa usoara, suportul mare pentru biblioteci si frameworkuri de deep learning, atat ca documentatie, cat si ca implementari efective.

De asemenea componenta de "model manager" care ruleaza intr-un thread fiecare model este scrisa in python.

5.5 Java

Java este un limbaj de programare orientat-obiect, puternic tipizat, conceput de catre James Gosling la Sun Microsystems (acum filial Oracle) la inceputul anilor 90, fiind lansat in 1995. Cele mai multe aplicatii distribuite sunt scrise in Java, iar noile evoluii tehnologice permit utilizarea sa si pe dispozitive mobile gen telefon, agenda electronic, palmtop etc. In felul acesta se creeaza o platforma unic, la nivelul programatorului, deasupra unui mediu eterogen extrem de diversificat. Acesta este utilizat in prezent cu succes si pentru programarea aplicatiilor destinate intranet-urilor

In proiect, Java a fost folosit la crearea backendului pentru interfata web. Motivatia alegerii este experienta puternica anterioara a autorului.

5.6 Spring Boot

In cadrul proiectului am pus la dispozitie si o interfata web in care se afiseaza reprezentarea vizuala a retelei generate de modelul ales de utilizator. Backendul acestei interfete este scris in Java si are in spate Spring Boot.

Pe scurt Spring Boot este frameworkul web pe care l-am folosit pentru a scrie serverul proiectului. Ni se pune la dispozitie posibilitatea de a scrie relativ usor o aplicatie web pe arhitectura MVC.

Motivul pentru care s-a ales Spring Boot in locul unui alt framework web de Java este documentatia vasta si experienta anterioara a autorului cu aceasta tehnologie.

5.7 RaphaelJS

RaphaelJS reprezinta biblioteca de javascript folosita pentru desenarea in browser retelei generate de modelul ales de utilizator. RaphaelJS este un

wrapper peste WebGL cu un API relativ usor de folosit.

5.8 Twitter Bootstrap

Twitter Bootstrap este un framework web pentru componenta de frontend. Se pun la dispozitie mai multe clase CSS pentru crearea de componente vizuale precum butoane, bare de meniu, cat si un mecanism specific de de pozitionare a continutului.

Motivul pentru care am introdus si aceasta dependinta in stiva de tehnologii folosita in proiect si nu am creat toata interfata in RaphaelJS este performanta. Desenarea intregului design si implementarea si mecanicilor din spate (precum gestionarea apasarilor de buton) ar fi fost remarcabil mai lenta.

6 Concluzii

În momentul de față, cele două modele antrenate, nu produc rezultate multumitoare, nu sunt capabile să genereze comparatori pentru rețele de dimensiuni mari. Optimizarea acestor modele pentru inputuri $N > 4$ poate reprezenta o direcție de viitor pentru proiect. Acest lucru ar putea fi realizat printr-o metodă introdusă relativ recent (*Mai 2018*)[18].

Readucem în atenție că această lucrare conține un element de noutate, întrucât această abordare de a genera rețele de sortare folosind *deeplearning* nu a mai fost atinsă până în acest moment.

Bibliografie

- [1] <https://github.com/primaryobjects/nnsorting>
- [2] [https : //github.com/drforester/Sequence_to_Sequence_Sorting](https://github.com/drforester/Sequence_to_Sequence_Sorting)
- [3] <https://github.com/aditya-prasad/dnnet>
- [4] <https://adventuresinevolutionblog.wordpress.com/2016/09/17/minimal-sorting-networks/>
- [5] <http://aclweb.org/anthology/I17-3017>
- [6] <http://psycnet.apa.org/record/1999-02657-007> – Seven times seven is about fifty
- [7] <https://www.sciencedirect.com/science/article/pii/S0020025512007670>
- [8] <https://github.com/apache/incubator-mxnet/tree/master/example/bi-lstm-sort> – sort numbers using lstm architecture
- [9] [https : //rylanschaeffer.github.io/content/research/neural_turing_machine/main.html](https://rylanschaeffer.github.io/content/research/neural_turing_machine/main.html)
- [10] <https://github.com/carpedm20/NTM-tensorflow>
- [11] <http://www.robots.ox.ac.uk/~tvlg/publications/talks/NeuralTuringMachines.pdf>
- [12] <https://medium.com/snips-ai/ntm-lasagne-a-library-for-neural-turing-machines-in-lasagne-2cdce6837315>
- [13] <http://www.iti.fh-flensburg.de/lang/algorithmen/sortieren/networks/nulleinsen.htm>
- [14] <http://www.cs.tau.ac.il/~zwick/Adv-Alg-2015/Sorting-Networks.pdf>
- [15] Valsalam, Vinod K. and Miikkulainen, Risto Using Symmetry and Evolutionary Search to Minimize Sorting Networks *J. Mach. Learn. Res.*, 14(1): 303–331, feb 2013
- [16] D. E. Knuth. Art of Computer Programming: Sorting and Searching, volumul 3, capitolul 5, paginile 219229. Addison-Wesley Professional, 2 edition, April 1998

- [17] M. W. Green. Some improvements in non-adaptive sorting algorithms. In Proceedings of the Sixth Annual Princeton Conference on Information Sciences and Systems, pagine 387-391, 1972.
- [18] Stephen McAleer, Forest Agostinelli, Alexander Shmakov, Pierre Baldi Solving the Rubik's Cube Without Human Knowledge, 18 Mai 2018
- [19] Alex Graves, Greg Wayne, Ivo Danihelka Neural Turing Machines, 20 Oct 2014
- [20] Peter Kipfer, Rdiger Westermann Improved GPU Sorting, *https :
//developer.nvidia.com/gpugems/GPUGems2/gpugems2_chapter46.html*,
Aprilie 2005