

Laborator 3: Retele Convolutionale - Convolutional Neural Networks (CNNs)

1. Introducere

Laboratorul trecut s-a studiat problema clasificarii unei baze de date simple (MNIST) cu ajutorul unei retele de tip perceptron multistrat (Multilayer Perceptron - MLP). Pe parcursul acestui studiu de caz s-au abordat pasii elementari pentru rezolvarea unei astfel de probleme (alegera arhitecturii, a functiei loss, a optimizatorului, a ratelor de invatare, impartirea bazei de date in batch-uri precum si bucla de invatare).

In continuare, solutia usoara din laboratorul trecut va fi extinsa prin folosirea arhitecturilor mult mai puternice de tip retele convolutionale (CNN). De asemenea se vor discuta si alte aspecte utile in utilizarea retelelor neuronale.

2. Motivatie si aspecte generale

Reteaua neuronală folosita anterior a utilizat drept trasaturi de intrare chiar pixelii constituenti ai imaginilor. Desi rezultatele au fost satisfacatoare, acest dezvoltamant nu ar fi fost la fel de probabil pentru baze de date mai complicate. In general, valorile directe ale pixelilor nu sunt considerate a fi trasaturi puternice. Cu acest scop se utilizeaza extractoare de trasaturi, precum HoG si LBP, care surprind mai bine informatia spatiala din fiecare zona de interes sau din jurul fiecarui pixel.

Retelele convolutionale au adus o serie de imbunatatiri in ceea ce priveste algoritmi de *machine learning*. O parte dintre acestea vor fi discutate ulterior. Drept punct de plecare se va face referire la schema generala a unei retele convolutionale:

Se pot observa doua zone principale:

- Prima, formata din straturile de tip Convolution si Pooling
- A doua, formata din straturi *Fully Connected*

Prima zona s-a dovedit a avea rolul de extractor de trasaturi. Primele straturi convolutionale extrag informatii de tip contururi, iar acestea devin mai complexe odata cu parcurgerea retelei. S-a observat ca trasaturile ieșite dupa prima zona a retelei sunt puternice (desi nu s-a descoptor ce inseamna), in general fiind mai reprezentative decat trasaturile obtinute prin aplicarea algoritmilor HoG sau LBP.

A doua zona, cu straturile *Fully Connected* este practic o retea MLP aplicata trasaturilor extrase de zona convolutionala a retelei.

2.1. Straturile convolutionale

Acest tip de strat este unitatea de baza din noile arhitecturi. Facand o analogie cu procedeele consacrate din prelucrarea imaginilor, un strat convolutional realizeaza o serie de operatiuni de filtrare liniara pe matricea de la intrarea in strat (pe rand). Aceasta matrice poate fi imaginea in sine, sau rezultatul altui strat, denumit *feature map* (harta de trasaturi). Un fapt relevant in acest caz este ca nucleul de convolutie are adancimea matricei de la intrare (ex. 3 pentru o imagine RGB). De la modul in sine in care functioneaza stratul convolutional se pot observa niste diferente fata de modul de functionare al MLP, precum si proprietati relevante:

- *sparse interactions* (interactiuni "rare"): spri deosebire de MLP, unde fiecare neuron din stratul N interactuia cu toti neuronii din stratul $N + 1$, la un CNR doar o serie de neuronii din stratal curent vor participa la neuronul din stratal urmator (zona denumita campul receptiv, *receptive field*). Acest fapt implica stocarea a mai putini parametri, si a operatiilor mai putine;
- partajarea parametrilor: nucleul de convolutie este folosit pentru intreaga imagine, deci ponderile care conduc la un anume neuron din stratal urmator sunt mereu aceleasi. Aceasta trasatura este pusă in contrast cu arhitectura MLP, unde fiecare pondere era folosita o singura data, pentru o perche anume de neuron;
- echivalarea a translate a reprezentarilor: se refera la proprietatea ca daca imaginea de intrare sufera o transformare de tip translate, si harta de trasaturi de iesire va suferi aceeasi modificare.

Toate straturile convolutionale sunt urmate de o functie de activare (in general Rectified Linear Unit - ReLU). Pentru a crea un strat convolutional in PyTorch se folosesc sintaxa:

```
conv_x = nn.Conv2d(in_channels = nr_canale_input, out_channels = nr_canale_output, kernel_size = [linii_filtru,coloane_filtru], stride = [pas_orizontal,pas_vertical], padding = [bordare_lini, bordare_coloane])
```

Forma corecta pentru datele intrare este de tipul [nr_imagine, canale, linii, coloane].

Atentie: cel mai frecvent, forma unei imagini obisnuite pentru alte biblioteci (scikit-image, OpenCV) este [linii,coloane,canale].

Argumentul *stride* se refera la pasul pe care il face nucleul convolutional dupa ce opereaza asupra zonei curente. Un pas de [1,1] inseamna ca va parcurge toti neuronii posibili. Argumentul *padding* se leaga de capetele imaginii (bordare). Se decide care este numarul de linii/coloane cu valori de zero care trebuie adaugate la imagine inaintea operatiei de convolutie. Aceasta operatie este folosita pentru a influenta dimensiunea matricei de iesire din stratal de convolutie.



Comparatie intre rezultatele unei operatii de convolutie, cu pas orizontal si vertical 1.

Se vede ca in varianta cu padding [1,1] (jos), iesirea are aceiasi dimensiune ca intrarea

2.2. Straturile de pooling

Operatia de *pooling* este cel mai bine tradusa in limba romana ca o "grupare", o subselecție. Pe scurt, aceasta operatie inlocuiește valoarea unei zone din imagine/harta de trasaturi cu o statistica a aceliei zone. Functia de *max-pooling* este cea mai folosita in aplicatii si inlocuiește valoarea dintr-o zona bine definita cu maximul aceliei zone, rezultand intr-o harta de trasaturi mai mica, dar care pastreaza cea mai relevanta statistica. In acest mod, pe langa reducerea dimensionalitatii, se obtine si o invarianta la translatiile mici.

Modul de creare a unui strat de *max-pooling* care sa ia vecinatia 2x2:

```
pool_x = nn.MaxPool2d(kernel_size = [linii_vecinatate, coloane_vechinatate], stride = [pas_orizontal, pas_vertical])
```



Rezultatul operatiei de *Max-pooling* cu argumentul *kernel_size* de valoare [2,2] si *stride* de valoare [2,2]

2.3. Straturile fully connected

Dupa cum a fost mentionat anterior, aceste straturi sunt cele *dens connectate*, obisnuite dintr-un MLP. Inainte de a le putea folosi, trebuie de trasaturi trebue aplatizate (vectorizate), folosind fie o metoda dedicata de "aplatizare":

```
flat = torch.flatten(resultat_strat_anterior, 1,3) # Se aplatizeaza dimensiunile 1-3 (adica se obtine ceva de dimensiunea canale x linii x coloane)
fie generand o varianta aplatizata folosind metoda implicita "view": flat = resultat_strat_anterior.view(-1,nr_canale*nr_lini*nr_coloane)
```

Atentie: ambele variante trebuie apelate direct in metoda *forward*, fiindca nu sunt straturi separate.

Pentru a crea un strat dens:

```
fc = nn.Linear(in_features = nr_neuroni_intrare, out_features = nr_neuroni_iesire)
```

Nu uitati ca, in general, pe ultimul strat nu se aplica functie de activare (de tip ReLU), iar *CrossEntropyLoss* aplica direct Softmax pe iesirea retelei.

3. Arhitecturi de baza

Există o multitudine de arhitecturi actuale, unele care au adus imbunatatiri marginale, altele care sunt specializate pe o anumita sarcina, dar cateva arhitecturi sunt considerate a fi pietre de temelie pentru domeniul. In continuare vor fi prezentate cateva arhitecturi care au atras o atentie foarte mare la momentul aparitiei lor.

3.1. LeNet-5

Cea mai veche arhitectura convolutionala a fost prezentata in 1998 cu scopul de a recunoaste cifre scrise de mana in documente. A fost conceputa pentru imagini de rezolutie mica (32 x 32 pixeli) si, din cauza constrangerilor vremii in ceea ce priveste puterea de calcul, nu a prezentat o adancime mare (doar 2 straturi convolutionale cu nuclee de 5 x 5 pixeli). Schema arhitecturii:

Arhitectura LeNet-5. Figura prelata din [1]

3.2. AlexNet

Dupa mai bine de un deceniu (in 2012), arhitectura AlexNet a fost prima arhitectura neuronală care a castigat concursul ILSVRC cu o arhitectura avand 5 straturi convolutionale, care nu era deosebit de considerabil mai mare, nuclee de convolutie mai mari in straturile initiale (11 x 11) cu pasi mai mari de parcurgere a imaginii si a folosit activari de tip ReLU. Aceasta arhitectura mult mai puternica decat ce-sa vehiculeaza pana in acel moment a fost antrenata cu ajutorul a doua GPU-uri performante. Schema arhitecturii:

Arhitectura AlexNet. Figura prelata din [2]

3.3. VGG

Urmatorul pas important adus in domeniu a fost studiul impactului adancimii unei retele. In acest scop, retelele din familia VGG au demonstrat cresterea performantei odata cu adancimea. Rețeaua VGG-19 (de la cele 19 straturi neuronale de o rînd) este printre cele mai mari retele utilizate in termeni de numar de parametri care trebuie invatati si cei intinuti minte in etapa de inferenta (propagare inainte a informatiei). In prezent, aceasta arhitectura este adevarata folosita pentru trasaturile generale puternice extrase dupa ultimul strat convolutional, utila si in alte sarcini decat clasificarea. O alta observatie importanta este reprezentata de reducerea tuturor nucleelor de convolutie la dimensiunea 3 x 3 cu pas 1 la deplasare. Schema arhitecturii VGG-16:

Arhitectura VGG-16. Figura prelata din [3]

3.4. ResNet

Aparuta mai recent ca celelalte arhitecturi prezentate, importanta acestei arhitecturi a fost uriasa, rezolvand problema *vanishing gradient*. Aceasta reprezinta scaderea puternica a gradientilor odata cu avansarea in reteaua in etapa de propagare inapoi. Practic, retelele cu un numar mare de straturi erau foarte greu de antrenat. Solutia acestui tip de arhitectura a fost introducerea blocului "residual", care presupunea ca la iesirea dintr-un bloc compus din mai multe convolutii se adauga si intrarea in bloc. Aceste conexiuni de tip "scurtcircuit" (sau "scurtcircuit") au permis crearea unei retele mult mai adanci (inclusiv 1000 de straturi). Ca o consecinta a numarului crescut de straturi, s-a putut scadea numarul de filtre per strat, pastrand numarul de parametri care trebuiau invatati relativ redus. Arhitectura ResNet-34 este ilustrata alaturi de VGG-19 (care are un numar considerabil mai mare de parametri):

Arhitectura ResNet-34 in comparatie cu VGG-19 si o arhitectura cu 34 de straturi fara "scurtcircuit". Figura prelata din [4]

Exercitiu. Incarcati baza de date MNIST si antrenati o retea convolutionala care are urmatoarea forma: convolutie + pooling + convolutie + pooling + 2 fully connected. Functiile modificate pentru citirea bazei de date se afla mai jos.

Atentie! Nu uitati ca forma ceruta de straturile de convolutie si tineti cont de consumul crescut de resurse la o retea de dimensiuni mai mari.

```
import numpy as np
# Nu trebuie sa te ocupi cu citirea datelor, dar trebuie pentru tot restul
import torch
def get_MNIST_train():
    mnist_train_data = np.zeros((16000,784))
```

```
    mnist_train_labels = np.zeros(16000)
```

```
    f = open("train-images.idx3-ubyte",'rb',encoding = 'latin-1')
```

```
    g = open("train-labels.idx1-ubyte",'rb',encoding = 'latin-1')
```

```
    byte = f.read(16) # bytes magic number, 4 bytes nr_imagine, 4 bytes nr linii, 4 bytes nr coloane
```

```
    byte_label = g.read(8) #bytes magic number, 4 bytes nr labels
```

```
    mnist_train_data = np.fromfile(f,dtype=np.uint8).reshape(16000,784)
```

```
    mnist_train_data = mnist_train_data.astype(np.float32)
```

```
    mnist_train_labels = np.fromfile(g,dtype=np.uint8)
```

```
    # Conversie pentru a se potriviri cu procesul de antrenare
    mnist_train_data = mnist_train_data.astype(np.float32)
    mnist_train_labels = mnist_train_labels.astype(np.int64)
```

```
    return mnist_train_data, mnist_train_labels
```

```
def get_MNIST_test():
    mnist_test_data = np.zeros((1000,784))
```

```
    mnist_test_labels = np.zeros(1000)
```

```
    f = open("t10k-images.idx3-ubyte",'rb',encoding = 'latin-1')
```

```
    g = open("t10k-labels.idx1-ubyte",'rb',encoding = 'latin-1')
```

```
    byte = f.read(16) # 4 bytes magic number, 4 bytes nr_imagine, 4 bytes nr linii, 4 bytes nr coloane
```

```
    byte_label = g.read(8) #bytes magic number, 4 bytes nr labels
```

```
    mnist_test_data = np.fromfile(f,dtype=np.uint8).reshape(1000,784)
```

```
    mnist_test_data = mnist_test_data.astype(np.float32)
```

```
    mnist_test_labels = np.fromfile(g,dtype=np.uint8)
```

```
    # Conversie pentru a se potriviri cu procesul de antrenare
    mnist_test_data = mnist_test_data.astype(np.float32)
    mnist_test_labels = mnist_test_labels.astype(np.int64)
```

```
    return mnist_test_data, mnist_test_labels
```