

ML – Brain scan detection

Polifronie Dragos-Constantin, 344

Rezolvarea acestui task se bazeaza classifierul multi-layer de perceptroni din biblioteca sklearn.

Au fost conduse mai multe teste in ceea ce priveste gasirea unei variante mai bune, printre care:

Modificarea numarului de layere, modificarea solverului si rata de invatare initiala, dar si prelucrarea in prealabil a datelor (in prima instanta normalizarea lor, urmand mai apoi testarea acuratetii dupa reducerea zgomotului din imagini si aplicare threshold – testate in late submission)

General:

- Imaginile au fost incarcate cu ajutorul bibliotecii opencv, iar procesul a fost grabit folosind multiprocessing.
- Imaginile au fost mai apoi normalizate printr-un StandardScaler din biblioteca sklearn, in functia `normalize_data(train_data, test_data, type='standard')`
- Intrucat imaginile au fost numeroase, iar timpul si resursele nu mi-au permis folosirea tuturor, doar un procent (ales convenabil) din pozele de training au fost folosite, fapt ce probabil a tras in jos performanta modelului. In submisile trimise, doar 10% din imagini au fost folosite pentru antrenare.

Parametrii testati inainte de submisii: (testing score se refera la pozele carora le stim clasa, adica poze din validation si /sau training dar care nu au fost folosite pentru antrenare, intrucat pozele de testare nu sunt clasificate)

In urmatoarele teste vor fi puse niste imagini ce redau rezultatele testarii mai multor functii de activare, impreuna cu solverul folosit ca parametru

- Primul test a fost facut pe un clasificator cu 2 layere a cate 100 de neuroni fiecare si cu functia de activare = identity.

Intrucat functia de activare identity este una lineara iar datele nu sunt, asa cum ne-am astepta, rezultatul nu este satisfactor, iar aceasta metoda nu este intocmai eficienta pentru modelul nostru.

Rezultate pentru functia de activare identity, pe o retea cu 2 layere a cate 100 neuroni.

Solver	Iterations	Trainin g loss	Testin g score
sgd - constant	200	nan	0.3873 52
sgd - invscaling	200	nan	0.3873 52
adam - init=0.01	57	0.2427 94	0.6482 21
adam - init=0.001	118	0.0395 42	0.6837 94
lbfgs - alpha=0.0001	200	0.1538 43	0.6521 74
lbfgs - alpha=0.00001	200	0.1162 40	0.6284 58

Se poate observa ca functia sgd nu a putut converge din cauza linearitatii functiei de activare.

- Rezultate pentru functia de activare tanh

Solver	Iterations	Trainin g loss	Testin g score
sgd - constant	200	0.3960 53	0.6798 42
sgd - invscaling	68	0.6578 34	0.5928 85
adam - init=0.01	30	0.6960 97	0.6126 48
adam - init=0.001	16	0.6740 90	0.4505 93
lbfgs - alpha=0.0001	200	0.0630 70	0.7351 78
lbfgs - alpha=0.00001	200	0.0735 99	0.6917 00

- Rezultate pentru functia de activare logistic:

Solver	Iterations	Trainin g loss	Testing score
sgd - constant	200	0.6533 97	0.7272 73
sgd - invscaling	14	0.7148 27	0.3873 52
adam - init=0.01	16	0.6883 89	0.6600 79
adam - init=0.001	37	0.6361 58	0.7233 20
lbfgs - alpha=0.0001	200	0.3521 43	0.6679 84
lbfgs - alpha=0.00001	200	0.3694 88	0.7035 57

- REZULTATE PENTRU FUNCTIA DE ACTIVARE RELU:

Solver	Iterations	Trainin g loss	Testing score
sgd - constant	200	0.1568 05	0.6917 00
sgd - invscaling	146	0.4308 30	0.7008 03
adam - init=0.01	67	0.5234 27	0.6877 47
adam - init=0.001	101	0.0020 98	0.7075 10
lbfgs - alpha=0.0001	200	0.0000 15	0.7351 78
lbfgs - alpha=0.00001	200	0.0000 08	0.7430 83

Observam ca in cazul acestei functii, rezultatele deja se imbunatatesc, asa ca functia relu va fi folosita pentru urmatoarele testari.

In urmatoarele teste, vom modifica numarul de neuroni din layere, folosind de aceasta data 3 layere a cate 200. (functia de activare ramane aceeaasi – relu)

Rezultate pentru activation = relu, hidden_layer_sizes = (200, 200, 200)

Solver	Iterations	Trainin g loss	Testing score
sgd - constant	200	0.0180 41	0.6956 52
sgd - invscaling	200	0.7693 14	0.5928 85
adam - init=0.01	69	0.3334 50	0.6996 05
adam - init=0.001	96	0.0009 65	0.7312 25
lbfgs - alpha=0.0001	115	0.0000 24	0.7628 46
lbfgs - alpha=0.00001	116	0.0000 09	0.7667 98

In afara de solverul sgd- invscaling, toate scorurile de testare au crescut. Vom incerca in urmatorul test o retea cu 3 layere a cate 800 de neuroni fiecare

Solver	Iterations	Trainin g loss	Testing score
sgd - constant	111	0.0081 59	0.6205 53
sgd - invscaling	200	0.7003 78	0.6086 96
adam - init=0.01	70	0.6407 21	0.6837 94
adam - init=0.001	53	0.0011 42	0.7430 83
lbfgs - alpha=0.0001	91	0.0000 66	0.6798 42
lbfgs - alpha=0.00001	90	0.0000 11	0.6877 47

Un ultim test, cu 3 layere, avand 800, 400, respectiv 200 neuroni:

Solver	Iterations	Trainin g loss	Testing score
sgd - constant	191	0.0079 07	0.7075 10
sgd - invscaling	200	0.7703 76	0.6047 43
adam - init=0.01	94	0.0527 71	0.6956 52
adam - init=0.001	112	0.0017 18	0.7470 36
lbfgs - alpha=0.0001	76	0.0000 63	0.6877 47
lbfgs - alpha=0.00001	78	0.0000 22	0.7193 68

Printre aceste rezultate, solverul adam cu rata initiala de invatare 0.001 a adus in medie cele mai bune rezultate, asa ca acesta a fost ales pentru submitii.

Avand in vedere aceste rezultate:

Prima submitie (cea prezentata in **model1.py**) care a primit scor **0.27** , a folosit 3 layere a cate 200 de neuroni, fara early stopping cu un tol de 0.001.

A doua submitie (cea prezentata in **model2.py**) a primit un scor de aprox **0.32** si a folosit 3 layere a cate 800 de neuroni, fara early stop, cu aceeasi valoare de tol.

A fost “testata” problema unui overfit, prin a 3a submitie care a folosit un early stopping = True, dar care a avut un rezultat mult sub primele, si anume 0.04

In **late submissions** (folosite pentru adunarea de informatii pentru documentatie), algoritmi au fost alterati, marind numarul de imagini folosite pentru antrenare si batchul fiecarei epoci in clasificatorul mlp.

De altfel, imaginile au fost preprocesate pentru a obtine date mai utile, prin reducerea zgomotului din imagini cu ajutorul unui filtru gaussian cu kernel 5x5 (a fost folosita biblioteca opencv). In aceste cazuri, totusi, acuratetea nu a fost imbunatatita, ci dimpotriva, diminuata, cele mai bune rezultate locale fiind de 0.38 iar publice 0.30

Imaginile au fost trecute printr-un blur gaussian iar mai apoi un threshold adaptiv a fost aplicat pentru a evidentia partile de interes din imagine. In acest caz, insa, poate fi vorba de un overfit, intrucat pe multimea de validare, acuratetea a ajuns chiar pana la 90%, dar cu toate acestea, la testele finale (submitii) aceste modele au obtinut un scor mai mic decat submitia nr 2.

In concluzie, modelul mlp propus ar putea fi imbunatatit pentru a aduce o acuratete mai mare prin prelucrarea mai buna a imaginilor (probabil un sharpening al imaginii inainte de threshold ar ajuta) si printr-o mai buna gestionare a evolutiei modelului, intrucat este posibil ca modelul sa fi ajuns la overfit si sa nu mai poata clasifica imaginile de testare. Cu toate acestea, modelul nu pare promitator in a atinge o acuratete foarte mare (majoritatea rezultatelor, indiferent de schimbarile aduse, se invart in jurul acuratetei de 0.3-0.4)

-biblioteci folosite:

- * Tensorflow 2.11.0
- * opencv 4.5.4
- * sklearn 1.0.2