

# Google Landmark Recognition Challenge

Lupașcu Marian

Coordonator: Iulia Duță

## Introducere

În primăvara lui 2018 Google a lansat Landmark Recognition Challenge, o competiție ce presupune antrenarea unui clasicator pentru recunoașterea clădirilor și monumentelor celebre de pe mapamond. În acest proiect am implementat o serie de modele printre care: VGG-16[1] preantrenat pe ImageNet la care am antrenat doar clasicatorul de la final și VGG-16 preantrenat modificat la care am adăugat un modul de atenție.

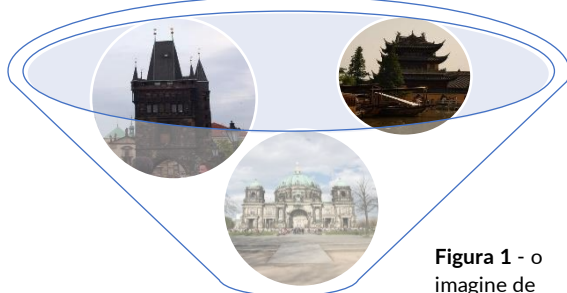


Figura 1 - o imagine de ansamblu asupra clasicatorului

Catedrala din Berlin + Jin Mao Building + Podul Carol

## Arhitectura rețelei

Pe date s-a făcut **augmentare**, anume: pe train: rotație la  $\pm 45^\circ$ , random crop la  $224 \times 224$ , schimbare de iluminare de până la 25% și random flip, pe validare și test: center crop la  $224 \times 224$  și rotație la  $\pm 45^\circ$ . Arhitectura folosită este VGG-16 la care s-a modificat secvența de clasificare de la final, folosind **Dropout** după fiecare layer de fully-connected pentru **regularizare**. Antrenarea modelului se face doar pe secvența de **clasificare**, secvența de **feature extraction** rămânând înghețată pe parcursul antrenării, deoarece VGG-16 fiind deja preantrenat pe 1000 de clase aste de așteptat ca aceasta parte să returneze o serie de feature-uri bine definite, urmând să fie clasificate.

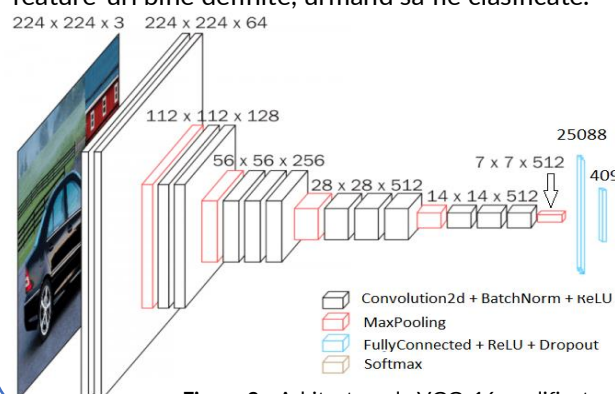


Figura 3 - Arhitectura de VGG-16 modificată

## Descrierea Datasetului

Datasetul este format din aproximativ 12 milioane de imagini cu 14952 clădiri (clase), pe care le-am descărcat de internet apoi am luat cele mai definite 27 de clase, adică aproximativ 11.5% din dimensiunea datasetului (46GB) și 0.18% din numărul de clase. O histogramă care prezintă greutatea celor mai definite 100 de clase din dataset este descrisă. Se poate observa cu ușurință că datele sunt dezechilibrate variind de la 49091 exemple pentru clasa 9633 până la 1 pentru câteva zeci de clase (9936, 99, etc.). Datasetul a fost împărțit în 3 părți: train(70%), validare(15%) și test(15%).

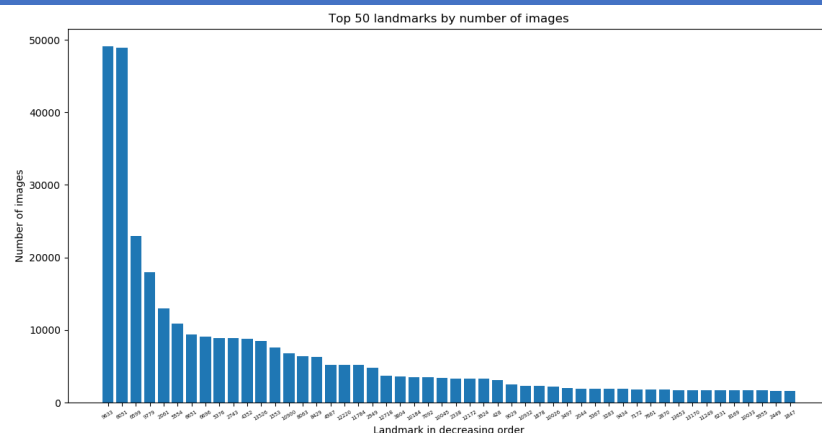


Figura 2 - Histograma în care este prezentat numărul de imagini din primele 50 de clase sortate după greutate.

## Modulul de atenție

- În continuare extindem rețeaua anterioară la un **Attention Branch Network (ABN)**[2] ceea ce nu este altceva decât VGG-16 cu structura de **branch** și un modul de atenție. Arhitectura de tip branch este sugerată abia după aplicarea funcției ReLU de după cele două convoluționale  $7 \times 7 \times 512$  cu filtru  $(3 \times 3)$ . Scopul **modulului de atenție** este de a forța rețeaua să se uite numai pe anumite zone din imagine. Menționez că arhitectura folosită este una similară cu cea descrisă în paperul - **Attention Branch Network: Learning of Attention Mechanism for Visual Explanation**[2].
  - Attention map**-ul se obține pe baza a patru convoluții din **feature map**, reducând numărul de canale de la 512 la numărul de clase (la noi 27) apoi prin aplicarea **sigmoidei**. Acesta este o hartă  $7 \times 7$  în cazul nostru de valori între 0 și 1, în care zonele cu scor aproape de 1 prezintă interes și cele aproape de 0 nu.
  - Prin înmulțirea fiecărui canal din **feature map** cu **attention map**-ul se obține un nou **feature map**.
- Loss**-ul este calculat ca o **combinație liniară** dintre partea de jos (**attention map** + clasificare) și partea de sus (**GAP** (funcționează pe ideea de **estompare** a informației)).

$$Loss_{final} = \alpha \cdot Loss_{Attention\ branch} + \beta \cdot Loss_{perception\ branch}$$

- $\alpha$  și  $\beta$  sunt **hiperparametrii** care măsoară cât valorează **atenția** (partea de sus) și cât **percepția** (partea de jos).

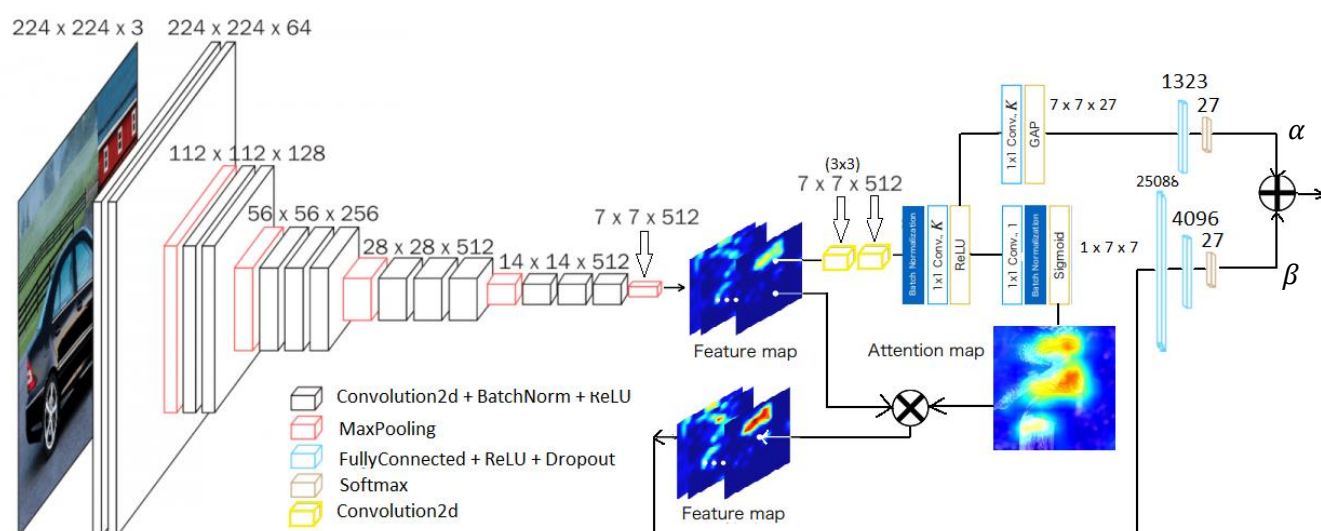


Figura 4 - Arhitectura de VGG-16 modificată peste care s-a adăugat modulul de atenție (overview asupra Attention Branch Network) [2]

## Experimente și Rezultate

Modelul	Validare	Test
VGG-16	87.92%	87.97%
VGG-16 cu modul de atenție $\alpha = 0.25, \beta = 1$	30.54%	27.66%
VGG-16 cu modul de atenție $\alpha = 0, \beta = 1$	32.71%	32.58%

Figura 5 - Acuratetea modelelor după 3 epoci pentru diferiți hiperparametrii

- În Figura 5 se observă diferența de acuratetea după 3 epoci între cele 3 modele, un motiv ar fi numărul mai mare de parametrii al ultimelor două modele.
- Un aspect interesant este dat de ultimele două linii din table care sugerează faptul că rețeaua care are **loss** format numai de pe ramura de clasificare (cea de percepție) are rata de convergență puțin mai mare comparativ cu rețeaua care a acumulat **loss** din ambele părți, sugerând astfel că informația estompata (generală) nu aduce un plus în cazul nostru de clasificare.

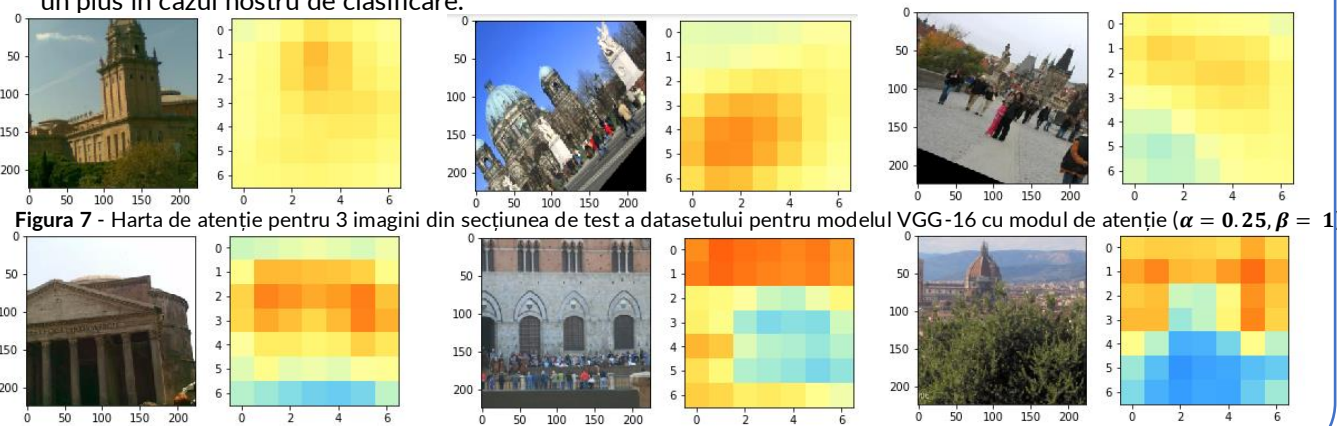


Figura 7 - Harta de atenție pentru 3 imagini din secțiunea de test a datasetului pentru modelul VGG-16 cu modul de atenție ( $\alpha = 0.25, \beta = 1$ )

Figura 8 - Harta de atenție pentru 3 imagini din secțiunea de test a datasetului pentru modelul VGG-16 cu modul de atenție ( $\alpha = 0, \beta = 1$ )

## Probleme la implementare

- Una din marile impedimente de care am dat la implementare a fost eroarea **CUDA out of memory**, atunci când trecem de la o epocă la alta. Soluție: antrenez doar câte o epocă apoi salvez modelul pe drive. La sfârșit de epocă încarc noii parametrii și repet procedul până termin numărul de epoci.
- Un alt impediment este durata mare de train, aproximativ 60-65 minute în acest moment. Menționez că a fost redus de la aproximativ 3h.

## Concluzii

- Atenția ajunge să se concentreze pe elementele de interes din imagine, însă învățarea este îngreunată de modelul mai complex și necesită o tunare mai atentă a hiperparametriilor și o atenție mai mare la procesul de optimizare.
- Modulul de atenție este un concept general care poate fi aplicat la orice task ce presupune modele convoluționale și poate fi aplicat la diverse nivele din rețea, aspect ce nu a apucat să fie tratat în această lucrare. (fiind pus la final de feature extraction)

## Referințe

- [1] - Attention Branch Network: Learning of Attention Mechanism for Visual Explanation - Hiroshi Fukui, Tsubasa Hirakawa, Takayoshi Yamashita, Hironobu Fujiyoshi
- [2] - Very Deep Convolutional Networks for Large-Scale Image Recognition - Karen Simonyan, Andrew Zisserman