

În aceste două săptămâni ne-am focusat pe studiul sistemelor de image recognition și tehnologiilor existente de image labeling / adnotation și am parcurs câteva tutoriale precum Google ML Crash Course și Tensor Flow Neural Network Course. Am parcurs diferite lucrări de cercetare legate de:

- Identificarea zonelor relevante din pagini web folosind DOM trees, însă acestea erau outdated, cele comprehensive fiind din ~2004.
- Annotated Transport Web Forms care crează un model pentru automatizarea cumpărării de bilete pentru transport.

Concluziile la care am ajuns sunt că pentru problema noastră nu există corpuri care să poată fi folosite în antrenarea modelului precum și faptul că ne-am putea folosi de service-ul de text recognition doar pentru partea de labeling a componentelor într-un stadiu mult mai avansat al proiectului.

În urma identificării problemei lipsei unui dataset pe care să lucrăm, soluția pe care am găsit-o e următoarea: facem un webcrawler care să salveze mai multe pagini cu caracteristici similare, apoi în paralel să realizăm

1. Clasificarea manuală a unor elemente de pe pagina (ex: partea de meniu, butonane, labeluri etc).
2. Dezvoltarea unui script: datele relevante de pe fiecare pagină vor fi extrase, apoi transformate în features relevante, ce pot fi folosite mai departe pentru antrenare. În legătură cu acest script am avea nevoie de sugestii: ne gândeam să dezvoltăm un script care să tag-uiască elemente folosind codul html din spate (sau chiar să angajăm indieni)

În săptămânile ce urmează vom încerca să ne împărțim următoarele task-uri:

- Dezvoltarea aplicației pentru adnotarea manuală a imaginilor de test și exportarea datelor în csv.
- Dezvoltarea unui webcrawler care să facă screenshot și să descarce codul html unui număr mare de pagini web. (scrapy + phantomJS)
- Dezvoltarea unui script care să automatizeze adnotarea elementelor folosind codul html.

În paralel cu crearea unui corpus vom încerca să dezvoltăm sistemul folosind CNN în Tensor Flow și să îl antrenăm pe un dataset deja existent (ex: autonomous driving ) pentru o înțelegere mai bună a funcționalității acestora.