

VYSOKÉ UČENÍ TECHNICKÉ V BRN

FAKULTA INFORMAČNÍCH TECHNOLOGIÍ

Projektová dokumentácia k predmetu SUR

Detektor osoby

25. apríla 2020

Dominik Bobo (xbobos00)
Tomáš Kender (xkende01)

1 Úvod

Zadaním projektu je natrénova detektor jednej osoby z obrázku tváre i hlasovej nahrávky. Detektor je implementovaný v programovacom jazyku Python 3.8.

2 Spustenie a používanie detektoru

2.1 Intalácia kniníc

Na spustenie vetkých súastí programu je potrebné ma nainstalované vetky potrebné kninice a súasti, ktoré sú obsiahnuté v súbore `requirements.txt`. Okrem toho je potrebné ma nainstalované aj linuxové kninice `sox` a `libsndfile1-dev`.

Intaláciu pythonových kniníc je možné uskutočniť pomocou `pip`:

```
$ pip install -r requirements.txt
```

Prípadne:

```
$ pip3 install -r requirements.txt
```

Poznámka: Program sa dá spustiť aj priamo cez priložený `dockerfile`, ktorý nainštaluje linuxové i pythonové balíky do virtuálneho prostredia a tým nerozhodí existujúcu konfiguráciu na počítači.

2.2 Dátová truktúra

Program oakáva nasledujúcu dátovú truktúru.

```
-- |-- data                // dáta určené k~trénovaniu a testovaniu detektoru
    |-- target_dev
    |-- non_target_dev
    |-- target_train
    |-- non_target_train
    |-- eval                // dáta určené ku klasifikácií
    |-- src                 // zdrojové kódy
requirements.txt
```

2.3 Spustenie jednotlivých súastí detektoru

Spustenie detektoru je možné pomocou príkazu.

```
$ python src/detector.py
```

Prípadne:

```
$ python3 src/detector.py
```

Po spustení sa natrénujú modely tréningovými dátami, modely sa ukladajú do priečinka `src` a následne sa uskutoční klasifikácia dát v `eval` a ukladajú sa výsledky do `.txt` súboru. Funkcie detektoru je možné upravovať pomocou prepínačov na príkazovom riadku.

```
--verbose    // informácie pre ladenie
--plot       // vykreslí informácie získané počas tréningovania
--evalonly   // vyhodnotia sa len dáta v~/eval (potrebné ma uložené modely)
--trainonly  // tréningovanie modelu bez vyhodnotenia
--cnn        // použitie CNN modelu pre spracovanie obrázkov
--system     // výber systému, ktorý sa má spustiť (face alebo voice)
```

3 Pouité modely

3.1 Modely pre spracovanie obrázkov

Pred samotným tréňovaním modelov je dôležité získať o najviac tréňovacích a testovacích dát. Toto je docieľené tým, že na tréňovanie sa používajú vetky *.png* súbory v priežku *data*. Obrázky sa načítajú pomocou knižnice *Pillow*. Target dátam sa nastaví label 1, ostatným label 0. Pomocou *ImageDataGenerator* z knižnice *Keras* sa vytvoria ďalšie tréňovacie dáta a to vďaka modifikáciám ako posun, rotácia, skrivenie, prevrátenie. Takto získané dáta smerujú modelom na tréňovanie.

3.1.1 Model CNN

CNN je model konvolučných neuronových sietí, bol vybraný na základe zvedavosti, či CNN dokáže na akceptovateľnej úrovni rozpoznať tváre. Tento model získava features obrázkov z rôznych vrstiev, známe aj ako *feature maps*, pomocou *pooling-u* postupne matice „zmenujeme“, respektíve podvzorkujeme. Nakoniec použijeme funkciu *flatten* čím vytvoríme 1D pole. Ako optimalizátor je použitý **SGD** - *Stochastic-Gradient Descent*, a aby sme zabránili pretréňovaniu, dáta sa predkladajú v náhodnom poradí a zo vetkých dát sa náhodne vyberie 20%, ktoré sa budú používať ako testovacie dáta. Tréňovanie modelu potom prebieha v desiatkach *Epochs* a samotné tréňovanie trvá pomerne dlho. Model sa po tréňovaní uloží `src/modelCNN.h5`.

3.1.2 Lineárny model SVM s SGD

Ako druhý model bol zvolený **SGDclassifier** použitý z knižnice *sklearn.linear_model*. Používa lineárny klasifikátor **SVM** - *Support Vector Machines* s tréňovaním dát SGD. Model bol vybraný hlavne ako o dos rýchlejšia alternatíva k CNN modelu vzhľadom na tréňovanie modelu a celková rýchlosť v porovnaní s ostatnými klasifikátormi a pritom so stále použitou úspešnosťou.

4 Spracovanie zvuku

4.1 Random Forest s MFCC

Z nahrávky sa najprv odstráni um, následne sa oddelí ticho. Pre extrahovanie features bola zvolená knižnica *librosa*. Z nej boli využit MFCC a *pitch* hlasu pre charakterizáciu nahrávky. MFCC featúry sa následne ešte alej kálovali a prehľadali funkciou *zscore* z knižnice *scipy*. O sa týka použitého klasifikátora, tu padla voľba na **Random Forrest Classifier** zo *sklearn.ensemble*. Jeho hlavnou výhodou je delegovanie úloh na jednotlivé rozhodovacie stromy. Každý rozhodovací strom spraví predikciu a najrozšírenejší výsledok sa zobere ako výsledok. Základným predpokladom ale je to, že jednotlivé stromy sú navzájom od seba nezávislé (ie používajú pri vyhodnocovaní rôzne vlastnosti), a teda aj keď niektoré stromy dojdú k zlému výsledku, prevažná väčšina dojde k správnejmu a túto chybu kvantitatívne vyváži. Z toho dôvodu platí, že čím vyší máme počet stromov, tým vyšiu úspešnosť máme na to, že prešíslime chybujúce stromy. Zároveň to ale pochopitene zvyšuje zaťaženie procesora. Druhou najvýznamnejšou konfigurovanou premennou je maximálna hĺbka stromu. Ladením a testovaním som nakoniec skončil pri lese s veľkosťou 96 stromov a vhodnou hĺbkou.

5 Záver

Pre detekciu tváre sa používajú celkovo dva modely. Presnosť CNN modelu je vyššia, čo je pravdepodobne spôsobené aj dlhším tréňovaním dát. Presnosť by sa dala zvýšiť zmenou určitých nastavení ako napríklad pridanie ďalšieho filtra alebo použitím vhodnejšieho optimalizátora. Pri SGD modeli je najväčším problémom nemožnosť nastaviť získanie score pre vyhodnotený obrázok, avak pre jeho rýchlosť a prijateľnú presnosť sme sa rozhodli ho ponechať v projekte. Z hľadiska zvýšenia presnosti by bolo najvhodnejšie tento model vynechať, nakoľko lineárny klasifikátor nemusí byť najvhodnejšia voľba pre naše dáta.