

Spis treści

1. Wstęp	3
2. Rozproszone uczenie na urządzeniu końcowym	4
2.1. Projekt systemu	5
2.1.1. Protokół	5
2.2. Algorytm optymalizacji	6
3. System weryfikacji użytkownika	8
3.1. Wstępne przetwarzanie obrazu	8
3.2. Weryfikacja użytkownika	8
3.2.1. Procedura weryfikacji	8
3.2.2. Ekstraktor cech	8
3.2.3. Zbiory danych	8
3.2.4. Metody weryfikacji twarzy	9
3.2.5. Rozporszona weryfikacja twarzy	10
Bibliografia	11
Wykaz symboli i skrótów	12

1. Wstęp

Coraz częściej urządzenia internetu rzeczy stają się głównymi urządzeniami komputerowymi coraz większej liczby użytkowników (cite). Często gromadzą one wrażliwe dane i dostęp do takich urządzeń przez niewłaściwe osoby grozi nieodwracalnymi stratami dla ich właściciela. Nowe przyrządy wyposażone w odpowiednie sensory pozwalają na uwierzytelnienie dostępu już nie tylko za pomocą hasła ale również przez weryfikację biometryczną. Zabezpieczenia biometryczne mogą się opierać na rozpoznawaniu linii papilarnych, głosu, skanowaniu żył, czy też tęczy lub siatkówki oka. W szczególności popularnym rozwiązaniem jest weryfikacja użytkownika przez biometrię twarzy (jakiś cytat).

Metody weryfikacji twarzy (ją jakieś) wymagają wyspecjalizowanych i dokładnych kamer (need fact check). Z jednej strony montowanie drogich i nowoczesnych kamer na urządzeniach IoT do celów weryfikacji jest nieopłacalne finansowo, a z drugiej z perspektywy użytkownika pożądanym jest posiadanie możliwie dokładnego systemu weryfikacji dostępu. Najnowocześniejsze i najbardziej dokładne metody bazują w całości lub przynajmniej części na bazie sieci neuronowych (cite,fact check). Metody te pewnym stopniu niwelują potrzebę posiadania wyspecjalizowanych kamer jednak dokładność nowoczesnych metod jest bardzo uzależniona od jakości i ilości danych, które posłużyły do wytrenowania sieci neuronowej.

Sensory, w które wyposażone są te urządzenia (na przykład aparat, mikrofon, itp), w połączeniu z faktem, że są używane codziennie, gromadzą niebywałą ilość, zazwyczaj prywatnych, danych. Modele wyuczone na takich danych dają znakomitą poprawę ich użyteczności jednak ze względu na wrażliwy charakter danych wiąże się z ryzykiem i wysoką odpowiedzialnością ich przechowywania w scentralizowanej lokalizacji albo nawet całkowitym brakiem dostępu do tych danych.

Federated Learning pozwala na bezpieczne dla użytkownika wykorzystanie jego prywatnych danych (zdjęć) w celu dotrenowania sieci neuronowych i w tym poprawy ich jakości. W tej pracy zostanie zbadana metoda ta metoda uczenia sieci neuronowych w implementacji systemu rozpoznawania twarzy systemu na urządzeniach IoT.

Główne kontrybucje tej pracy to 1) Implementacja i weryfikacja algorytmu Federated Averaging dla zadań klasyfikacji obrazów oraz weryfikacji twarzy.

2. Rozproszone uczenie na urządzeniu końcowym

W tym rozdziale przedstawiony zostanie projekt systemu, który korzystając z podejścia Federated Learning (FL) umożliwi na trening ekstraktora cech twarzy w rozproszonym środowisku urządzeń IoT.

Federated Learning Federated Learning (FL) jest podejściem zastosowania uczenia maszynowego w rozproszonym środowisku pozwalające na trening modeli na dużym zbiorze danych zdecentralizowanych prywatnych danych znajdujących się na urządzeniach końcowych użytkowników, a w szczególności urządzeniach Internetu Rzeczy. FL realizuje ideę "przyniesienie kodu do danych, zamiast danych do kodu" i adresuje fundamentalne problemy prywatności, własności i lokalności danych. Federated learning został opisany w (cite).

Problemy odpowiednie do zastosowania federated learningu mają następujące właściwości:

- 1) Trening na rzeczywistych danych gromadzonych na urządzeniach mobilnych daje znaczącą przewagę nad treningiem na ogólnie dostępnych danych proxy dostępnych w centrach danych.
- 2) Te dane są prywatne albo są zbyt duże do przetrzymywania ich w centrach danych
- 3) Dla zadań nadzorowanych, etykiety danych powstają samoistnie z interakcji użytkownika z urządzeniem.

Cechy środowiska systemu Algorytmy optymalizacji mogące być zastosowane do optymalizacji na urządzeniach IoT mają kilka cech wyróżniających je od znanych już algorytmów rozproszonej optymalizacji:

- **Non-IID** Dane trenujące na danym urządzeniu są zazwyczaj zależne od konkretnego użytkownika i dlatego lokalny zbiór danych zebrany na dowolnym urządzeniu nie będzie reprezentatywny w stosunku do dystrybucji całej populacji
- **Niezbalansowany** Podobnie, niektórzy użytkownicy będą o wiele częściej korzystali z aplikacji aparatu niż inni, co będzie prowadziło do różnic w wielkości zebranych lokalnych zbiorów danych trenujących.
- **Masywnie rozproszony** Spodziewa się, że liczba finalnych użytkowników biorąca udział w optymalizacji będzie większa niż średnia liczba przykładów trenujących przypadająca na jednego klienta.
- **Ograniczona komunikacja** Urządzenia IoT są pomimo założenia, że mają dostęp do internetu mogą być ograniczone wolnym albo kosztownym łączem sieciowym.

W tej pracy główna uwaga zostanie poświęcona na doprowadzenie systemu do działania w patologicznym przypadku środowiska danych Non-IID oraz ograniczonej i zawodnej komunikacji.

2.1. Projekt systemu

Implementowany system umożliwia trening głębokich sieci neuronowych na danych przetrzymywanych na bezpośrednio urządzeniu IoT. Dane te nigdy nie opuszczają samego urządzenia. Wagi modelu są agregowane i łączone na serwerze znajdującym się w chmurze za pomocą algorytmu Federated Averaging konstruując nowy model globalny, który zostaje przesłany z powrotem do urządzeń końcowych do inferencji. System ten został opisany pierwotnie w (cytat) i z sukcesem zaaplikowany w implementacji inteligentnej klawiatury smartphona. (cite blogaska googla).

Przedstawienie architektury systemu zaczynamy od przedstawienia protokołu według, którego przebiega cały przepływ danych. Opisujemy to w następnej sekcji.

2.1.1. Protokół

Uczestnikami protokołu są urządzenia Internetu Rzeczy oraz serwer, który jest usługą chmurową. Urządzenia anonsują serwerowi swoją gotowość do uruchomienia zadania. Zadanie jest specyficzne dla populacji urządzeń i polega na lokalnym uruchomieniu obliczeń, takich jak trening z zadanymi przez serwer parametrami albo ewaluacja wytrenowanego modelu na lokalnych dla urządzenia danych.

Protokół dzieli komunikację na rundy, z której każda zaczyna się od selekcji. Z potencjalnie tysięcy gotowych urządzeń w danym oknie czasowym, serwer wybiera ich podzbiór, który dostanie zaproszenia do wzięcia udziału w obliczeniach. Liczebność tego podzbioru jest parametrem algorytmu opisywanego w kolejnym rozdziale

Serwer mówi wybranym urządzeniom jakiego typu obliczenia będą wykonywane. Wraz z tą informacją przesłany jest graf obliczeniowy modelu, zserializowane parametry globalnego modelu oraz pozostałe informacje niezbędne do uruchomienia obliczeń w danej rundzie. Następnie każdy z uczestników wykonuje lokalnie obliczenia bazując na globalnym stanie oraz swoim lokalnym zbiorze danych i przesyła wynik obliczeń z powrotem na serwer. W szczególności wynikiem, w przypadku rundy treningowej, będą wytrenowane parametry modelu lokalnego. Serwer wciela zebrane aktualizacje w swój stan globalny co kończy pojedynczą rundę i proces ten się powtarza.

Opisywany protokół pozwala urządzeniom końcowym na poprawianie globalnego, pojedynczego modelu pomiędzy rundami, z których każda składa się z trzech faz co zostało pokazane na rysunku (rysunek?).

Selekcja Cyklicznie, urządzenia, które spełniają wymagane kryteria (Patrz sekcja XXX) meldują serwerowi swoją gotowość. Serwer wybiera podzbiór z dostępnych mu urządzeń. W naszej implementacji wybrano prostą, losową selekcję z rozkładem jednostajnym N urządzeń z tych, które ogłosiły swoją gotowość.

Konfiguracja Serwer jest skonfigurowany na podstawie wybranego mechanizmu selekcji urządzeń oraz agregacji modeli. Urządzenia końcowe są samo-konfigurowane według

Algorytm 1 FederatedAveraging. The K clients are indexed by k ; B is the local minibatch size, E is the number of local epochs, and η is the learning rate.

Server executes:

```

initialize  $w_0$ 
for each round  $t = 1, 2, \dots$  do
     $m \leftarrow \max(C \cdot K, 1)$ 
     $S_t \leftarrow$  (random set of  $m$  clients)
    for each client  $k \in S_t$  in parallel do
         $w_{t+1}^k \leftarrow \text{ClientUpdate}(k, w_t)$ 
     $w_{t+1} \leftarrow \sum_{k=1}^K \frac{n_k}{n} w_{t+1}^k$ 

```

ClientUpdate (k, w): // Run on client k

```

 $\mathcal{B} \leftarrow$  (split  $\mathcal{D}_k$  into batches of size  $B$ )
for each local epoch  $i$  from 1 to  $E$  do
    for batch  $b \in \mathcal{B}$  do
         $w \leftarrow w - \eta \nabla \ell(w; b)$ 
    return  $w$  to server

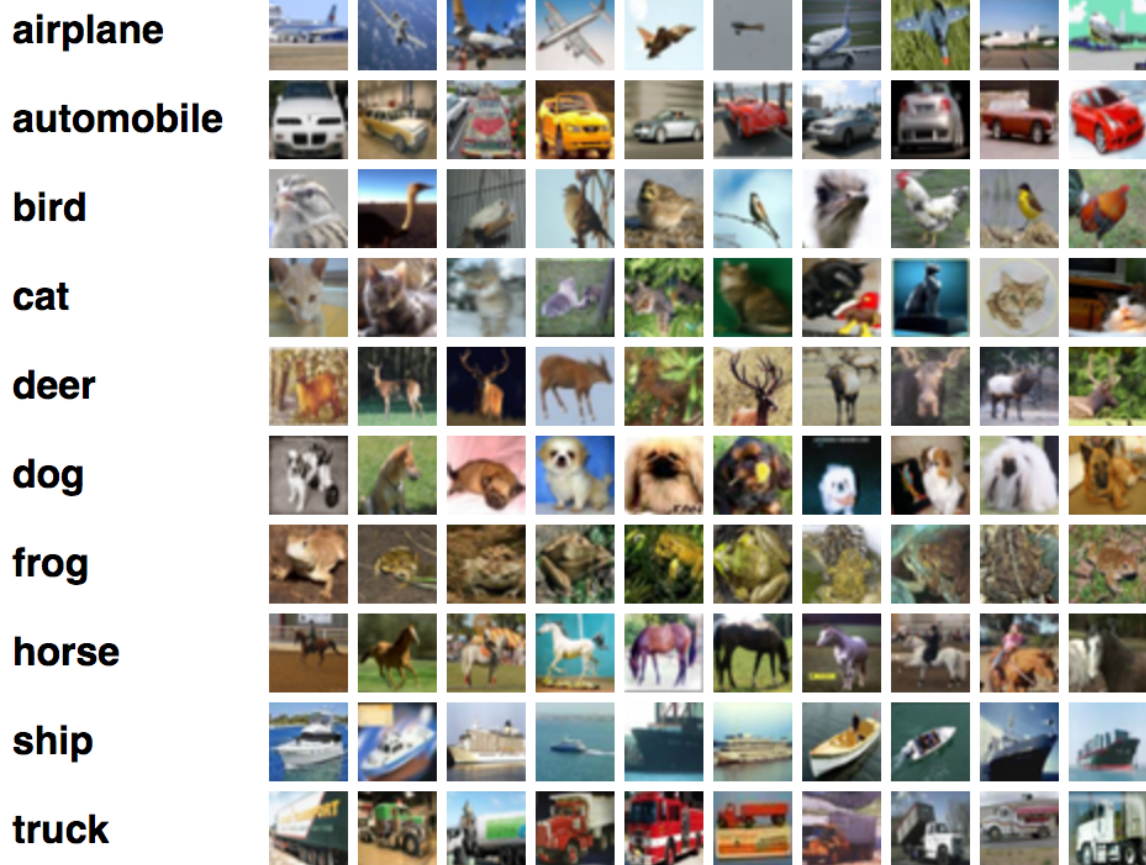
```

przesłanej przez serwer informacji o typie obliczeń, grafu obliczeniowego oraz stanu globalnego.

Raportowanie Serwer czeka, aż partycypujące urządzenia zdadzą raport z wynikiem przeprowadzonych obliczeń. Jeśli runda była rundą treningową to wraz z otrzymaniem wszystkich zaktualizowanych modeli, serwer agreguje je z użyciem algorytmu Federated Averaging i informuje raportujące urządzenia kiedy mogą ponownie zgłosić swoją gotowość. Jeśli wystarczająca liczba urządzeń zaraportuje wyniki serwerowi w pewnym oknie czasowym, runda zostanie uznana za zakończoną z powodzeniem i zagregowany model zastąpi aktualny model globalny. W przeciwnym przypadku zebrane wyniki z tej rundy zostają porzucone. Protokół jest w pewnym stopniu odporny na to, że partycypujące urządzenia z pewnych powodów nie odpowiedzą na przesłaną konfigurację albo nie zaraportują wyników.

Fazy selekcji i raportowania są dospecyfikowane przez zestaw parametrów. Podczas fazy selekcji serwer bierze pod uwagę pożądaną liczbę uczestników rundy, czasy time-outów oraz minimalną liczbę potrzebnych uczestników do rozpoczęcia rundy. Faza selekcji trwa dopóty dopóki nie zostanie wybrana docelowa liczba urządzeń albo nie skończy się czas przeznaczony na fazę selekcji. W drugim przypadku runda zostanie rozpoczęta jeśli została zebrana minimalna liczba urządzeń, w przeciwnym przypadku runda zostaje porzucona. Faza raportowania parametryzowana jest przez minimalną liczbę urządzeń, które muszą zaraportować wynik.

2.2. Algorytm optymalizacji



Rysunek 1. 10 przykładowych obrazów dla każdej z 10 klas zbioru CIFAR10

Algorytm 1 opisany został zaimplementowany w języku Python. Do implementacji modeli neuronowych i algorytmów uczących został wykorzystany framework PyTorch [1]. Poprawność implementacji została sprawdzona na zadaniu klasyfikacji obrazów wykorzystując prostą sieć konwolucyjną oraz popularny zbiór danych CIFAR10.

Sieć konwolucyjna Jako obiekt treningu omawianego algorytmu zastała użyta niewielka sieć neuronowa zawierająca dwie warstwy konwolucyjną z filtrami o szerokości 5x5 (pierwsza z 32 kanałami, druga z 64, po każdej dodatkowa warstwa 2x2 max pooling), po których następuje dwu-warstwowy perceptron i na końcu warstwa przekształcenia liniowego, co daje w sumie 10^6 miliona parametrów. Model został podsumowany w tabeli

CIFAR10 CIFAR10 jest popularnym syntetycznym zbiorem danych. Zbiór danych składa się z 60 000 kolorowych obrazów podzielonych na 10 klas, z 6000 obrazami przypadającymi na jedną klasę. Zawarte są w nim obrazy o szerokości i wysokości 32 pikseli. Standardowo zbiór dzieli się na dwa zbalansowane klasowo podzbiory: testowy i treningowy zawierających odpowiednio 10000 i 50000 etykietowanych przykładów. Na rysunku 1 znajduje się 10 losowo wybranych obrazów, dla każdej z 10 klas.

Protokół treningowy Do sprawdzenia poprawności implementacji została zaimplementowana procedura treningowa wzorowana na [2]. Zbiór treningowy został podzielony pomiędzy 100 użytkowników tak żeby każdy zawierał po 500 przykładów trenujących. Z powodu braku naturalnego podziału danych na tak dużą liczbę klientów rozważany jest tutaj nieco mniej wymagający przypadek, w którym dane każdego użytkownika są zbalansowane oraz równomiernie rozdyskrebowane.

Naszym celem była maksymalizacja dokładności z jaką model klasyfikował obrazy pochodzące ze zbioru testowego. Badanie jakości końcowego modelu globalnego odbywało się już nie w sposób rozproszony, a na serwerze stosując cały dostępny zbiór testowy, co zostało umożliwione dzięki wprowadzonym uproszczeniom co do rozkładu danych.

Obrazy uległy standardowemu przetworzeniu wstępnemu i augmentacji. Przykłady trenujące zredukowano do wielkości 24x24 pikseli przez losowe obcięcie krawędzi, obrazy uległy losowemu horyzontalnemu odbiciu lustrzanemu oraz standardowej normalizacji.

Zaimplementowany algorytm został porównany do standardowego algorytmu SGD (cytat).

Ewaluacja

3. System weryfikacji użytkownika

Zadaniem systemu jest weryfikacja użytkownika, czy zostały przedstawione w następnej sekcji.

3.1. Wstępne przetwarzanie obrazu

3.2. Weryfikacja użytkownika

Weryfikacja twarzy jest zadaniem przyrównania twarzy kandydata to innej, i weryfikacja czy nastąpiło ich dopasowanie. Jest to mapowanie jeden-do-jednego: należy sprawdzić czy jest to ta sama osoba.

3.2.1. Procedura weryfikacji

3.2.2. Ekstraktor cech

tutaj opisz jaki jest pipeline wchodzi zdjęcie i wychodzi embedding

3.2.3. Zbiory danych

Są jakieś zbiory danych 1) VggFace2 2) MS1M 3) MegaFace 4) LWF

Sprawdzić wpływ datasetu początkowego na wyniki ewaluacji początkowego modelu globalnego

ms1m do pre treningowania ponieważ posiada dużą różnorodność twarzy oraz ogólną dużą liczbę przykładów trenujących. VggFace2 zostanie z kolei wykorzystany do końcowego

dotreniowania modelu - ze względu na dużą liczbę zdjęć przypadająca na jedną osobę posiada dużą różnorodność wewnątrzklasową co powinno zwiększyć finałową efektywność modelu - większa odporność na rotację twarzy względem obiektywu, zmienne warunki świetlne itp.

Ze względu na popularność dwóch ostatnich zbiorów w ewaluacji systemów rozpoznających twarze zostaną one właśnie wykorzystane do ewaluacji, MegaFace ma tę dodatkową zaletę, dużej liczby osób, a dwa pierwsze tylko do treningu.

3.2.4. Metody weryfikacji twarzy

Wyjściem z modelu jest wektor(embedding), który pozwala na odróżnienie jednej twarzy od drugiej za pomocą porównania embeddingów. Papier google'a twierdzi, że udało im się uzyskać dobry score, jednak nie rzadkie inne prace nie raportują sukcesywnego otrzymania podobnego wyniku. Może to wynikać z 1) braku zbliżonego wielkością zbioru danych użytych do uzyskania wyniku raportowanego przez google'a 2) brakiem podobnych zasobów obliczeniowych jednej twarzy do embeddingu drugiej twarzy. Coś o tym, że są tak jakby dwie rodziny algorytmów. Jedna bazująca na klasycznym podjęciu stosowanym podczas klasyfikacji obrazów i druga bazująca na optymalizacji multi-class classification z użyciem cross-entropy loss (jak to do uja wafla przetłumaczyć?).

Papier google'a twierdzi, że udało im się uzyskać dobry score, jednak nie rzadkie inne prace nie raportują sukcesywnego otrzymania podobnego wyniku. Może to wynikać z 1) braku zbliżonego wielkością zbioru danych użytych do uzyskania wyniku raportowanego przez google'a 2) brakiem podobnych zasobów obliczeniowych, jednak nie rzadkie inne prace nie raportują sukcesywnego otrzymania podobnego wyniku. Może to wynikać z 1) braku zbliżonego wielkością zbioru danych użytych do uzyskania wyniku raportowanego przez google'a 2) brakiem podobnych zasobów obliczeniowych.

Triplet loss jest o wiele lepszy dla zastosowań federated learningu procedura trenowania

Classification loss Pokazać, że jest kilka rodzajów takich lossów 1) Cross entropy loss 2) Cosine loss 3) Arc loss

Implementacja Papier google'a twierdzi, że udało im się uzyskać dobry score, jednak rzadkie inne prace nie raportują sukcesywnego otrzymania podobnego wyniku. Może to wynikać z 1) braku zbliżonego wielkością zbioru danych użytych do uzyskania wyniku raportowanego przez google'a 2) brakiem podobnych zasobów obliczeniowych

Została podjęta próba implementacji wytrenowania modelu stosując w.w.y metodę wykorzystując zbiór danych VggFace2.

W celach pokazania, że implementacja jest poprawna porównujemy następujące modele

3. System weryfikacji użytkownika

Porównanie: dwa pojęcia: dwa datasety, modele niekoniecznie moje ale testowane na własnej ewaluacji

triplet(google report), wytrenowanych na danych prywatnych triplet(own), centre(external source) wytrenowanych na VggFace2 triplet(own), centre(external source) wytrenowanych na MS1m -> wyjdzie, że triplet loss jest o wiele mniej skuteczna

see this: <https://arxiv.org/pdf/1901.08616.pdf> <https://arxiv.org/pdf/1709.02940.pdf>

3.2.5. Rozporszona weryfikacja twarzy

porównanie:

centre loss ms1m

globalny model arcface uczony na ms1m, testowany na lwf i megaface globalny model arcface uczony na ms1m, dotrenowany tripletem na VggFace2, testowany

globalny model arcface dotrenowanie na lwf ale negatywne przykłady samplujemy z generatora globalny model arcface dotrenowanie na megaface ale negatywne przykłady samplujemy z generatora

rozporoszonej optymalizacji:

- **Non-IID** Dane trenujące on danym urządzeniu są zazwyczaj zależne od konkretnego użytkownika i dlatego lokalny zbiór danych zebrany na dowolnym urządzeniu nie będzie reprezentatywny w stosunku do dystrybucji całej populacji
- **Niezbalansowany** Podobnie, niektórzy użytkownicy będą o wiele częściej korzystali z aplikacji aparatu niż inni, co będzie prowadziło do różnic w wielkości zebranych lokalnych zbiorów danych trenujących.
- **Masywnie rozporszony** Spodziewa się, że liczba finalnych użytkowników biorąca udział w optymalizacji będzie większa niż średnia liczba przykładów trenujących przypadająca na jednego klienta.
- **Ograniczona komunikacja** Urządzenia IoT są pomimo założenia, że mają dostęp do internetu mogą być ograniczone wolnym albo kosztownym łączem sieciowym.

W tej pracy główna uwaga zostanie poświęcona na rozwiązanie doprowadzenie systemu do działania w środowisku danych Non-IID oraz ograniczonej komunikacji.

Bibliografia

- [1] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga i A. Lerer, **Automatic Differentiation in PyTorch**, w *NeurIPS Autodiff Workshop*, 2017.
- [2] H. B. McMahan, E. Moore, D. Ramage, S. Hampson i B. A. y Arcas, *Communication-Efficient Learning of Deep Networks from Decentralized Data*, 2016. arXiv: 1602. 05629 [cs . LG] .

Wykaz symboli i skrótów

EiTI – Wydział Elektroniki i Technik Informacyjnych

PW – Politechnika Warszawska