

WSI LISTA 4

Bohdan Tkachenko 256630

June 13, 2023

1 Zadanie 1: Las decyzyjny dla bazy danych MNIST

Celem tego zadania było stworzenie i wytrenowanie lasu decyzyjnego do rozpoznawania ręcznie pisanych cyfr z bazy danych MNIST za pomocą biblioteki TensorFlow.

1.1 Metodologia

Dane MNIST zostały wczytane i przygotowane do analizy. Każdy obrazek (28x28 pikseli) został przekształcony do jednowymiarowej tablicy o długości 784. Następnie stworzono i wytrenowano las decyzyjny za pomocą funkcji RandomForestClassifier z biblioteki scikit-learn. Las składał się z 100 drzew decyzyjnych.

1.2 Wyniki

Dokładność modelu została oceniona na zbiorze testowym. Dokładność, mierzona jako procent poprawnie sklasyfikowanych obrazków, wyniosła 0.9705. Oznacza to, że las decyzyjny poprawnie rozpoznał 97.05% cyfr w zbiorze testowym.

1.3 Wnioski

Wyniki pokazują, że las decyzyjny może być skutecznym narzędziem do rozpoznawania ręcznie pisanych cyfr. Jednakże, wyniki te powinny być porównane z wynikami innych modeli, takich jak sieci neuronowe, aby dokładnie ocenić ich skuteczność.

2 Zadanie 2: Klasteryzacja zbioru danych MNIST za pomocą algorytmu DBSCAN

Celem tego zadania było zastosowanie algorytmu DBSCAN do klasteryzacji zbioru danych MNIST, który zawiera obrazy ręcznie pisanych cyfr. DBSCAN

został wybrany ze względu na jego zdolność do identyfikacji klastrów o dowolnym kształcie i rozmiarze, co jest szczególnie przydatne w przypadku danych MNIST, gdzie różne style pisanie mogą prowadzić do znacznej zmienności w obrębie tej samej klasy.

Dane MNIST zostały najpierw przekształcone do jednowymiarowych wektorów, a następnie znormalizowane. Następnie zastosowano redukcję wymiarowości za pomocą PCA, aby zmniejszyć wymiarowość danych do 50 głównych składowych.

Parametry DBSCAN, takie jak `eps` (maksymalna odległość między dwoma próbkami, aby były uważane za sąsiadów) i `min_samples` (minimalna liczba próbek w sąsiedztwie, aby punkt danych był uważany za punkt centralny), zostały dobrane eksperymentalnie. Ostatecznie, `eps` zostało ustawione na 10.0, a `min_samples` na 5.

Algorytm DBSCAN zidentyfikował 30 klastrów, co jest równo z maksymalną liczbą 30 sugerowaną w zadaniu. Zidentyfikowano również 10976 punktów szumu, co stanowi mniejszą część danych w porównaniu do poprzednich prób.

Dokładność klasyfikacji, obliczona jako procent obrazów, które zostały przypisane do klastra reprezentującego ich prawdziwą klasę, wyniosła 0.16. Jest to poprawa w porównaniu do poprzednich prób, ale nadal jest to dość niska wartość, co sugeruje, że wiele obrazów zostało niepoprawnie sklasyfikowanych.

Wnioski: DBSCAN, mimo że jest potężnym algorytmem klasteryzacji, może nie być najlepszym wyborem dla danych o wysokiej wymiarowości, takich jak MNIST. Wyniki mogą być poprawione przez dalsze dostrojenie parametrów DBSCAN, zastosowanie innej metody redukcji wymiarowości, lub zastosowanie innego algorytmu klasteryzacji.