

Detekcja ogórków grunotwych w środowisku naturalnym przy użyciu algorytmu YOLO

Jan Tyc

April 2023

1 Wstęp

Pomimo ogromnego udziału rolnictwa w gospodarce światowej oraz ciągłego zastępowania pracy ludzkiej przez maszyny, wciąż wiele prac związanych z uprawami jest wykonywanych ręcznie. Wiele owoców i warzyw jest zbieranych przez maszyny rolnicze, które zapewniają dużą wydajność, szybkość i oszczędności zarówno w małych jak i dużych gospodarstwach. Ze względu na specyfikę upraw wiele gatunków wciąż zbieranych jest ręcznie z uwagi na brak odpowiedniej metody, która zebrałaby plony nie uszkadzając ich oraz nie pozostawiłaby zniszczeń na polu. Dzięki gwałtownemu rozwojowi sztucznej inteligencji coraz częściej podejmowane są próby stworzenia maszyny do automatycznego i precyzyjnego zbioru, a w ciągu ostatnich kilku lat wydajność tych systemów uległa znacznej poprawie [3]. Oczywistymi wyzwaniami przy opracowywaniu takiego systemu wykrywania upraw są niespójne oświetlenie, nierozróżnialne tło, zachodzące na siebie lub zasłonięte plony, czy niskie rozdzielczości zdjęć na których przeprowadzana jest detekcja obiektów [4, 6, 7]. W poniższym raporcie postaramy się dokonać szerokiego przeglądu literatury w temacie optymalizacji algorytmu YOLO V7 oraz ogólnego problemu detekcji owoców i warzyw w środowisku naturalnym.

Dalsza część raportu jest ułożona w następujący sposób: rozdział 2 dotyczy objaśnienia działania rodziny algorytmów YOLO (na przykładzie YOLO V7), w rozdziale 3 zawarte są informacje na temat danych użytych do trenowania algorytmu, rozdział 4 dotyczy przeglądu literatury i potencjalnych rozwiązań wprowadzonych w algorytmach detekcji obiektów w celu ich optymalizacji, rozdział 5 to podsumowanie zgromadzonej wiedzy i prezentacja planu dalszej części badań.

2 Wstęp do YOLO

YOLO (You Only Look Once) to rodzina algorytmów detekcji obiektów, która charakteryzuje się obecnie najlepszymi metrykami dokładności wykrywania obiektów oraz bardzo dużą szybkością przetwarzania obrazu [5]. Najnowsze odsłony są

w stanie zwracać położenia obiektów na obrazie w tempie przekraczającym 100 FPS co sprawia, że algorytm może być wykorzystywany do przetwarzania danych w czasie rzeczywistym [1]. Algorytm składa się z trzech części: "Ogona" (eng. Backbone), "Szyi" (eng. Neck) oraz "Głowy" (eng. Head)

2.1 "Ogon"

Pierwszą częścią każdego algorytmu YOLO jest "Ogon". Jest to pierwsza warstwa do której wprowadzany jest obraz, który następnie przechodzi przez wiele warstw konwolucji, normalizacji batchy, funkcji aktywacji (Relu, Leaky Relu, Mish, Silu), Maxpoolingu czy konkatencji wcześniejszych warstw. W ten sposób z obrazu ekstrahowane są cechy, które następnie będą dalej przetwarzane i posłużą do detekcji obiektu. Na przestrzeni rozwoju algorytmu stosowano "ogony" o różnej architekturze (Darknet24, CSPDarknet53, Modified CSP v7 [5]) natomiast w wersji YOLO V7 zaproponowano sieć E-ELAN [1], która rozwiązuje problem zanikającej ścieżki gradientu, zachowując przy tym odpowiednią złożoność obliczeniową.

2.2 "Szyja"

Po wstępnym przetworzeniu dane przechodzą do "Szyi" gdzie będą dalej przetwarzane. Architektura "Szyi" różni się znacząco w zależności od wersji YOLO. W YOLO V7 dane przechodzą przez moduł SPPCSP oraz przez warstwy konwolucji i konkatencji wcześniejszych warstw.

2.3 "Głowa"

Algorytm działa w oparciu o następujące cztery podejścia:

1. Podział obrazu na sekcje
2. Rysowanie 'bounding boxes'
3. Wyliczenie IOU
4. Non-maximum-suppression

2.3.1 Podział na sekcje

Rozpoczynamy od podzielenia oryginalnego obrazu na komórki siatki $N \times N$ o jednakowym kształcie. Każda komórka siatki jest odpowiedzialna za zlokalizowanie i przewidzenie klasy obiektu, który obejmuje, wraz z wartością prawdopodobieństwa/pewności.

2.3.2 Rysowanie 'bounding boxes'

Kolejnym krokiem jest wyznaczenie bounding boxów, które odpowiadają prostokatom zaznaczającym wszystkie obiekty na obrazie. Możemy mieć tyle bounding boxów, ile jest obiektów w obrebie danego obrazu. YOLO określa atrybuty tych pól ograniczających za pomocą pojedynczego modułu regresji w następującym formacie, gdzie Y jest ostateczna reprezentacja wektorowa dla każdego pola ograniczającego. $Y = [bx, by, bh, bw, c1, c2, \dots, cn]$ gdzie:

- bx, by - koordynaty środka prostokata,
- bh, bw - wysokość i szerokość prostokata,
- $c1, \dots, cn$ - prawdopodobieństwo że w prostokacie znajduje się obiekt danej klasy

2.3.3 IOU

W większości przypadków, pojedynczy obiekt na obrazie może mieć wiele kandydatów do przewidywania, nawet jeśli nie wszystkie z nich są istotne. Celem IOU (wartość między 0 a 1) jest odrzucenie takich prostokątów, aby zachować tylko te, które są istotne.

2.3.4 Non-max-suppression

Ustawienie progu dla IOU nie zawsze jest wystarczające, ponieważ obiekt może mieć wiele pól z IOU poza progiem, a pozostawienie wszystkich tych pól może zawierać szum. Tutaj możemy użyć NMS, aby zachować tylko pudełka z najwyższym wynikiem prawdopodobieństwa wykrycia.

3 Dane użyte do treningu

Zbiór danych tworzą 594 zdjęcia ogórków gruntowych w środowisku naturalnym. Zdjęcia zostały zrobione w wielkopolsce w Lipcu 2022 r. na dwóch różnych polach. Do ich zrobienia użyto aparatu z telefonu Xiaomi Redmi Note 8 Pro, a rozdzielczość każdego zdjęcia to 4624x3472. Dane zostały oetykietowane tzn. na każdym zdjęciu zostały zaznaczone ogórki, które następnie posłużą jako 'ground truth' w uczeniu algorytmu.

4 Przegląd literatury

Odpowiednio szybka i dokładna detekcja owoców i warzyw w środowisku naturalnym jest kluczowym elementem w stworzeniu automatycznej i precyzyjnej maszyny do zbioru plonów. Problemy związane z detekcją obiektów możemy podzielić na dwie klasy:

- Związane ze zdjęciami na których będzie przeprowadzana detekcja (dane)

- Związane z algorytmem oraz jego treningiem

W dalszej części przyjrzymy się jakie największe wyzwania czekają nas w każdej z klas oraz jakie ich rozwiązania są obecnie dostępne.

4.1 Problemy z danymi potrzebnymi do detekcji

Algorytmy detekcji uczą się na zdjęciach na których zaznaczone są położenia obiektów, które wytrenowany algorytm powinien zlokalizować. W przypadku błędnego lub niedokładnego zaznaczenia danych treningowych algorytm będzie się źle uczył i nigdy nie osiągnie zadowalającej skuteczności detekcji, dlatego szczególnie ważnym jest dbanie o wysoką jakość zbioru uczącego. W przypadku małej ilości danych, mogą one zostać poddane augmentacji, czyli zmianie zdjęcia np. poprzez obrót, zmianę wielkości, odbicie lustrzane, zmianę saturacji i jasności zdjęcia, w celu zwiększenia ilości próbek uczenia [2].

4.2 Problemy związane z algorytmami detekcji obiektów

Powstanie maszyny do automatycznego zbioru jest dyktowane szybkością z jaką algorytm jest w stanie znajdować warzywa i owoce na zdjęciach z kamery. Jeśli jedno zdjęcie jest przetwarzane zbyt długo to maszyna nie jest wydajna, dlatego detekcja obiektów musi odbywać się w tempie przynajmniej 30-40 FPS, żeby móc konkurować z człowiekiem. Algorytm podczas uczenia może wpadać w minima lokalne, przez co może nie osiągać zadowalających rezultatów i tracić możliwość generalizacji, dlatego bardzo ważny jest odpowiedni dobór architektury i hiperparametrów uczenia. Sposób podawania zdjęć podczas treningu też ma duży wpływ na metryki uczenia. "Curriculum learning" to technika polegająca na stopniowym podawaniu do uczenia coraz to trudniejszych przypadków w celu powolnej i głębszej generalizacji algorytmu.

4.3 Inne podejścia problemu detekcji

Stosuje się również inne techniki detekcji owoców i warzyw, nieoparte na uczeniu głębokim [3]. Detekcja za pomocą koloru to podejście, które znalazło zastosowanie w detekcji warzyw i owoców o kontrastowym kolorze (w odniesieniu do otoczenia) i powstały już pierwsze maszyny do zbioru, które obecnie są w fazie testów [8]. Stosuje się również detekcje za pomocą tekstury oraz detekcje kształtu na obrazie [3].

5 Podsumowanie i plan dalszych prac i eksperymentów

Algorytm YOLO V7, zaprezentowany we wcześniejszej części raportu, jest obecnie najlepszym dostępnym algorytmem do detekcji obiektów, który oferuje najlepsze rezultaty dla szybkości od 5 do 160 FPS. Zmiana struktury warstw algorytmu, może okazać się nieskuteczna z uwagi na to, że architektura algorytmu została starannie dobrana przez twórców i wszelkie zmiany mogą nieznacznie

pogarszać metryki. Wobec tego najlepszym podejściem w przypadku optymalizacji detekcji może okazać się optymalizacja samego etapu treningu sieci. Dobranie odpowiednich hiperparametrów uczenia takich jak liczba epok, learning rate, decay itp. może znacząco przyczynić się do zwiększenia poprawności detekcji. W przypadku problemu detekcji małych obiektów, przetestowana zostanie metoda rozpoznawania małych ogórków na podstawie detekcji koloru żółtych kwiatów. Nienaznacznie zwiększając długość działania algorytmu, możemy znacznie poprawić metryki detekcji dla małych obiektów dla tego konkretnego przypadku. Z uwagi na małą liczbę danych zostanie przeprowadzona augmentacja danych w celu zwiększenia ich ilości i poprawienia generalizacji modelu. Żeby zapobiec problemowi ciężko-rozpoznawalnych przypadków zostanie przetestowana metoda oparta na Curriculum Learning tzn. na początku podawania do sieci łatwych przypadków detekcji (ogórek w całości widoczny, dobrze oświetlony), a następnie po wstępnym uczeniu, douczanie sieci trudniejszymi przypadkami (ogórki nachodzące na siebie, częściowo zasłonięte, z kiepskim oświetleniem) zapewni większą generalizację uczenia. Pomimo wielu trudności związanych ze stworzeniem maszyny do automatycznego zbioru owoców i warzyw, postęp technologiczny i osiągnięcia w dziedzinach detekcji obiektów i przetwarzania obrazów pozwalają zachować optymizm w dalszych pracach nad rozwojem tego rozwiązania. Stworzenie takiej maszyny może znacząco zmienić rynek rolniczy na świecie, zmniejszyć koszty ponoszone przez producentów żywności co może bezpośrednio przełożyć się na spadek cen żywności na świecie.

6 Bibliografia

References

- [1] YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors, Chien-Yao Wang and Alexey Bochkovskiy and Hong-Yuan Mark Liao, 2022, 2207.02696, arXiv, cs.CV
- [2] R. Fernández, H. Montes, J. Surdilovic, D. Surdilovic, P. Gonzalez-De-Santos and M. Armada, "Automatic Detection of Field-Grown Cucumbers for Robotic Harvesting," in IEEE Access, vol. 6, pp. 35512-35527, 2018, doi: 10.1109/ACCESS.2018.2851376.
- [3] Xiao, F.; Wang, H.; Li, Y.; Cao, Y.; Lv, X.; Xu, G. Object Detection and Recognition Techniques Based on Digital Image Processing and Traditional Machine Learning for Fruit and Vegetable Harvesting Robots: An Overview and Review, *Agronomy* 2023, 13,639, <https://doi.org/10.3390/agronomy13030639>
- [4] Gai Rongli, Chen Na, Yuan Hai, A detection algorithm for cherry fruits based on the improved YOLO-v4 model, 2021/05/26, Neural Computing and Applications, <https://doi.org/10.1007/s00521-021-06029-z>

- [5] Diwan, T., Anirudh, G. Tembhurne, J.V. Object detection using YOLO: challenges, architectural successors, datasets and applications. *Multimed Tools Appl* 82, 9243–9275 (2023). <https://doi.org/10.1007/s11042-022-13644-y>
- [6] Lawal, Mubashiru Olarewaju. "Tomato detection based on modified YOLOv3 framework." *Scientific Reports* 11.1 (2021): 1-11.
- [7] Liu, Guoxu, et al. "YOLO-tomato: A robust algorithm for tomato detection based on YOLOv3." *Sensors* 20.7, 2020
- [8] Muhammad Hammad Malik, Ting Zhang, Han Li, Man Zhang, Sana Shabbir, Ahmed Saeed, Mature Tomato Fruit Detection Algorithm Based on improved HSV and Watershed Algorithm, *IFAC-PapersOnLine*, Volume 51, Issue 17, 2018, Pages 431-436, ISSN 2405-8963, <https://doi.org/10.1016/j.ifacol.2018.08.183>.