A

Project Report

on

# DYNAMIC E-COMMERCE DATA EXTRACTION

Submitted in partial fulfilment of the requirements for the award of the degree of

**Bachelor of Technology**

in

**COMPUTER SCIENCE AND ENGINEERING**

by

**Bontha Laxmi Priya Reddy**
**(20EG105306)**

**Naghma Mulla**
**(20EG105332)**

**Thukkani Dinesh Reddy**
**(20EG105350)**

**Palreddy Sai Sriya Reddy**
**(20EG105722)**

Under the guidance of

**Mr. G. Kiran Kumar**
Assistant Professor

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**
**Venkatapur (V), Ghatkesar (M), Medchal (D), T.S - 500088**
**2023-24**

## DECLARATION

We hereby declare that the report entitled **"Dynamic E-commerce Data Extraction"** submitted to the **Anurag University** in partial fulfilment of the requirements for the award of the degree of **Bachelor  of Technology (B. Tech)** in Computer Science and Engineering is a record of  an original work done by us under the guidance of **Mr. G. Kiran Kumar, Assistant Professor** and this report has not been submitted to any other university for the award of any other degree or diploma.

<div align="right">

Bontha Laxmi Priya Reddy
20EG105306

Naghma Mulla
20EG105332

Thukkani Dinesh Reddy
20EG105350

Palreddy Sai Sriya Reddy
20EG105722

</div>

Place: Anurag University, Hyderabad
Date:

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING**

**CERTIFICATE**

This is to certify that the project entitled **"Dynamic E-commerce Data Extraction"** being submitted by **Bontha Laxmi Priya Reddy** bearing the Hall Ticket number **20EG105306, Naghma Mulla** bearing the Hall Ticket number **20EG105332, Thukkani Dinesh Reddy** bearing the Hall Ticket number **20EG105350** and **Palreddy Sai Sriya Reddy** bearing the Hall Ticket number **20EG105722** and in partial fulfilment of the requirements for the award of the degree of the **Bachelor of Technology** in **Computer Science and Engineering** in **Anurag University** is a record of bonafide work carried out by them under my guidance and supervision from academic year 2023 to 2024.

The results presented in this project have been verified and found to be satisfactory. The results embodied in this project report have not been submitted to any other University for the award of any other degree or diploma.

Internal Guide                                                                      Dean, CSE
Mr. G. Kiran Kumar                                                          Dr. G. Vishnu Murthy
Assistant Professor

External Examiner

# ACKNOWLEDGEMENT

We would like to express our sincere thanks and deep sense of gratitude to project supervisor **G. Kiran Kumar**, **Assistant Professor, Department of Computer Science and Engineering**, Anurag University for his constant encouragement and inspiring guidance without which this project could not have been completed. His critical reviews and constructive comments improved our grasp of the subject and steered to the fruitful completion of the work. His patience, guidance and encouragement made this project possible.

We would like to express special thanks to **Dr. V. Vijaya Kumar, Dean School of Engineering,** Anurag University, for his encouragement and timely support in our B. Tech program.

We would like to acknowledge our sincere gratitude for the support extended by **Dr. G. Vishnu Murthy**, **Dean, Department of Computer Science and Engineering,** Anurag University.

We also express our deep sense of gratitude to **Dr. V. V. S. S. S. Balaram**, Academic coordinator. **Dr. T. Shyam Prasad,** Assistant Professor**,** Project Coordinator and Project review committee members, whose research expertise and commitment to the highest standards continuously motivated us during the crucial stage of our project work.

Bontha Laxmi Priya Reddy
20EG105306

Naghma Mulla
20EG105332

Thukkani Dinesh Reddy
20EG105350

Palreddy Sai Sriya Reddy
20EG105722

# ABSTRACT

Dynamic e-commerce data extraction is pivotal in today's digital landscape, enabling businesses to swiftly gather real-time insights from diverse online sources to inform decision-making, market analysis, and competitive intelligence. This methodology involves the automated retrieval of data from dynamic web pages, where content is frequently updated or modified, posing unique challenges compared to static websites. Leveraging advanced web scraping techniques like DOM parsing, XPath, and CSS selectors, in conjunction with dynamic rendering technologies such as headless browsers or JavaScript execution engines, facilitates the extraction of desired data elements despite the dynamic nature of the underlying web content. Additionally, the utilization of APIs provided by e-commerce platforms or third-party data providers streamlines the extraction process by offering structured access to data repositories. The extracted data encompasses various facets of e-commerce operations, including product information like pricing, availability, and descriptions, alongside customer reviews, market trends, competitor pricing strategies, and sales analytics. However, challenges such as anti-scraping measures implemented by websites, data inconsistency stemming from website updates, and legal considerations regarding data usage and privacy regulations necessitate meticulous planning and adherence to ethical data collection practices. Despite these obstacles, dynamic e-commerce data extraction empowers businesses with invaluable insights to refine pricing strategies, enhance product offerings, boost customer engagement, and maintain competitiveness in the rapidly evolving digital marketplace. As e-commerce continues to expand and data-driven decision-making becomes increasingly indispensable, dynamic data extraction remains fundamental to successful e-commerce operations, fostering innovation and strategic growth on a global scale.

**Keywords:** Web Scraping, Machine Learning, Natural Language Processing (NLP), Data Extraction

V

# TABLE OF CONTENTS

# List of Figures

# List of Tables

| Table No. | Table Name | Page No. |
|:---:|:---|:---:|
| 2.1 | Comparision of Existing Method from selected Strategies | 7 |

**List of Abbreviations**

| Abbreviations | Full Form |
|---|---|
| DOM | Document Object Model |
| API | Application Program Interface |
| CNN | Convolutional Neural Network |
| NLP | Natural Language Processing |
| RegEX | Regular Expressions |
| UI | User Interface |
| UX | User Experience |
| YOLO | You Only Look Once |
| OCR | Optical Character Recognition |
| BERT | Bidirectional Encoder Representations from Transformers |
| EDA | Exploratory Data Analysis |
| IoU | Intersection over Union |
| AP | Average Precision |
| MAP | Mean Average Precision |

# 1.INTRODUCTION

In today's fast-paced digital marketplace, dynamic e-commerce data extraction plays a pivotal role in empowering businesses with real-time insights essential for informed decision-making and strategic planning. Unlike static websites, dynamic web pages undergo frequent updates and modifications, posing unique challenges for data extraction. Advanced web scraping techniques, including DOM parsing, XPath, and CSS selectors, coupled with dynamic rendering technologies such as headless browsers or JavaScript execution engines, enable the extraction of valuable data elements despite the dynamic nature of web content. Additionally, leveraging APIs provided by e-commerce platforms or third-party data providers streamlines the extraction process by offering structured access to data repositories. Extracted data encompasses crucial aspects of e-commerce operations, including product information, customer reviews, market trends, competitor strategies, and sales analytics. Integrating machine learning algorithms and natural language processing techniques enhances data extraction by automating categorization, sentiment analysis, and trend identification from unstructured textual data. However, challenges such as anti-scraping measures, data inconsistency, and legal considerations around data usage and privacy regulations underscore the importance of ethical data collection practices. Despite these challenges, dynamic e-commerce data extraction empowers businesses to optimize pricing strategies, enhance product offerings, improve customer engagement, and maintain competitiveness in the rapidly evolving digital landscape. As e-commerce continues to expand and data-driven decision-making becomes indispensable, dynamic data extraction remains a cornerstone of successful e-commerce operations, driving innovation and strategic growth in the global marketplace.

## 1.1. Motivation

Traditional e-commerce data collection methods are inefficient and error-prone, hindering businesses' competitiveness. To overcome these challenges, we've developed an innovative solution utilizing computer vision and natural language processing to extract data directly from e- commerce websites, revolutionizing analytics and empowering businesses with actionable insights. Our project aims to democratize access to e-commerce data, driving growth and innovation in the online commerce landscape.

## 1.2. Problem Definition

Acquiring accurate product details from online shopping sites is complex. Our project aims to create an efficient system using advanced technology for text comprehension and image recognition. This system will adapt to website changes, improving the online shopping experience by swiftly providing reliable information. Creating an efficient e-commerce data extraction system involves diverse components. Web scraping and API integration gather data, while text comprehension and image recognition process information. Adaptation strategies and reinforcement learning ensure system flexibility.

## 1.3. Objective Of The Project

The objective of the "Dynamic E-commerce Data Extraction" project is to automate the process of extracting data from e-commerce websites directly, enabling efficient access to product information. By leveraging web scraping and data extraction techniques, the project aims to gather comprehensive data sets from various e-commerce platforms. The extracted data includes product details such as pricing, descriptions, reviews, and availability. Furthermore, the project aims to provide analytics and insights derived from the extracted data, facilitating informed decision-making for businesses and consumers alike. By automating the data extraction process and delivering actionable insights, the project aims to streamline e-commerce operations, enhance market intelligence, and optimize the shopping experience for users.

## 1.4. Problem Illustration

The problem illustration for the project "Dynamic E-commerce Data Extraction" involves a comprehensive exploration of various models and algorithms utilized to address the challenge of extracting data from e-commerce websites. Traditional web scraping techniques typically involve manual coding to navigate through web pages and extract relevant information. While effective for simple tasks, these methods often lack scalability and robustness when dealing with dynamic and complex web structures.

Machine learning-based web data extraction tools leverage algorithms such as decision trees, random forests, and support vector machines to automate data extraction processes. These models are trained on labeled datasets to learn patterns and relationships within the web data, enabling them to accurately extract information from diverse sources. However, their performance may vary depending on the quality and diversity of the training data, as well as the complexity of the target websites.

| Model | Algorithm/Technique |
|---|---|
| R-CNN | Dataset: ImageNet<br>Classification: Binary SVM |
| Fast R-CNN | RoI pooling |
| Faster R-CNN | RPN |
| Mask R-CNN | RoI align (pixel level segmentation) |
| RetinaNet | ResNet, Feature Pyramid Network (FPN) |
| Single Shot Detector (SSD) | Single deep neural network, feed forward convolutional network |
| Histogram of Oriented Gradients (HOG) | Detection window, RoI |
| Region-Fully Convolutional Network (R-FCN) | Convolutional + RoI |
| Spatial Pyramid Pooling (SPP) | Pyramid pooling |
| Yolo | Classification and regression, Deformable Parts Model (DPM) and R-CNN |

Fig 1.1: Various models and algorithms used for e-commerce data extraction

In contrast, deep learning-based approaches, such as neural networks and Convolutional Neural Networks (CNNs), offer a more advanced and adaptive solution for web data extraction. These models can automatically learn hierarchical representations of web data, capturing intricate features and structures that may be challenging to extract using traditional methods. By leveraging techniques like Natural Language Processing (NLP) and image recognition, deep learning models can effectively parse text and extract information from product descriptions, reviews, and images on e-commerce websites.

The comparison of traditional, machine learning, and deep learning-based web data extraction tools involves evaluating their performance across various metrics, including accuracy, efficiency, and scalability. While traditional methods may suffice for simple extraction tasks, they are often labor-intensive and prone to errors. Machine learning-based approaches offer improved automation and accuracy but may require significant manual feature engineering and tuning. In contrast, deep learning-based approaches excel in handling complex and unstructured data but may require large amounts of labeled data and computational resources for training. By examining the

strengths and limitations of each approach, researchers and practitioners can determine the most suitable methodology for their specific e-commerce data extraction needs.

| Tool | Extraction rule | Technique | Precision (%) | Self-healing |
|---|---|---|---|---|
| TSIMMIS[29] | Wrapper-based | Traditional/statistical | Not available | No |
| WebOQL[30] | Tag tree | Traditional/statistical | Not available | No |
| WHISK[31] | Regular expression | Supervised learning | 69 | No |
| RAPIER[32] | Logic rules | Supervised learning | 89 | No |
| SRV[33] | Logic rules | Supervised learning | 58 | No |
| SoftMealy[34] | Regular expression | Supervised learning | 58 | No |
| DEPTA[35] | Tag tree | Un-supervised learning | 98 | No |
| Trinity[36] | Regular expression | Un-supervised learning | 96 | No |
| DeLA[37] | Regular expression | Un-supervised learning | 80 | No |
| OLERA[38] | Regular expression | Semi-supervised learning | 99 | No |
| Proposed system | Object detection | Deep learning | To be determined | Yes |

Fig 1.2: Comparison of traditional, machine learning, and deep learning-based web data extraction tools

# 2. <u>LITERATURE SURVEY</u>

The technologies highlighted in the literature survey represent diverse approaches to data extraction from online sources, each with its own set of advantages and challenges.

- **Regular Expressions (RegEx)**: Bernstein et al. (2003) advocate using RegEx for extracting specific data formats from well-structured websites. RegEx is a powerful tool for pattern matching, allowing for efficient extraction of targeted information. However, its brittleness arises from its reliance on fixed patterns, making it prone to failure when websites undergo layout changes.
- **Natural Language Processing (NLP)**: Ma et al. (2017) utilize NLP techniques such as Named Entity Recognition (NER) and Information Extraction (IE) to analyze product descriptions. NLP offers the advantage of adaptability to varying data structures and languages, enhancing accuracy and flexibility in data extraction. However, implementing NLP techniques requires a sophisticated infrastructure and consumes higher computational resources, posing scalability challenges.
- **Computer Vision and Image Recognition**: Zhi et al. (2019) employ computer vision and image recognition strategies to gain insights beyond textual data. This approach provides a comprehensive understanding of products by analyzing visual elements. However, challenges related to image quality, lighting variations, and complex backgrounds may affect the accuracy and reliability of the extracted data.
- **Web Scraping**: Smith et al. propose web scraping as a direct method for data collection from websites. Web scraping enables real-time updates and comprehensive information retrieval by accessing data directly from its source. However, it is susceptible to changes in website layout, which can disrupt the scraping process and lead to incomplete or inaccurate data extraction.
- **API Integration**: Patel and Lee suggest API integration for direct access to structured data. API integration offers efficient retrieval of information, bypassing the need for manual extraction. However, limitations in data access may arise depending on the availability of API endpoints, potentially restricting the scope of obtainable information.

In summary, each technology offers unique benefits and challenges in the context of data extraction from online sources. The selection of an appropriate strategy depends on factors such as the nature of the data, the complexity of the website structure, resource availability, and scalability requirements. The proposed strategies for data extraction from online sources present various advantages and limitations. Bernstein et al. (2003) advocate Regular Expressions for extracting specific data formats from well-structured websites efficiently, yet its brittleness makes it vulnerable to layout changes, necessitating frequent updates. Conversely, Ma et al. (2017) utilize NLP techniques like Named Entity Recognition and Information Extraction for more adaptable data extraction, though this approach requires a complex infrastructure and higher computational resources. Zhi et al. (2019) employs computer vision and image recognition for a comprehensive understanding of products beyond textual data, facing challenges related to accuracy influenced by image quality and environmental factors. Smith et al. propose web scraping for direct data collection, offering real-time updates but susceptible to layout changes, potentially affecting data accuracy. Lastly, Patel and Lee suggest API integration for direct access to structured data, providing efficient retrieval but may face limitations in data access depending on API endpoints' availability, potentially restricting information scope. Each approach offers unique benefits and challenges, highlighting the need for careful consideration of the context and requirements when selecting a data extraction strategy.

| Author(s) | Strategies | Advantages | Disadvantages |
|---|---|---|---|
| Mika A. Bernstein et al., 2003 | Utilizes Regular Expressions patterns to match specific data formats within HTML code. | Efficient for well-structured websites with consistent HTML coding. | Brittle and prone to failure with website layout changes. |
| Jin Ma et al., 2017 | NLP techniques like Named Entity Recognition (NER) and Information Extraction (IE). | More robust and adaptable to website changes compared to RegEx. | Requires more complex infrastructure and computational resources. |
| Chen Zhi et al., 2019 | Computer Vision and Image Recognition | Provides deeper product insights beyond textual data. | Accuracy can be impacted by image quality, lighting, and complex backgrounds. |
| Smith Etal | Web Scraping | Direct collection from websites | Suspectable to website layout changes |
| Patel & Lee | API Integration | Direct Access to Structure Data | Limited data access depending on API endpoints |

Table 2.1: Comparison of Existing Method from selected Strategies

# 3. <u>PROPOSED METHOD</u>

The proposed method for addressing the challenges in dynamic e-commerce data extraction involves a multi-faceted approach

## 3.1. WEBSITE SNAPSHOTS:

Website snapshots involve capturing visual representations of e-commerce websites to analyze layout structures and visual elements. This technology essentially takes a snapshot or screenshot of a web page at a particular moment in time, preserving its appearance and design. These snapshots serve as valuable resources for understanding the layout of a website, including the placement of various elements such as text, images, buttons, and navigation menus.

By analyzing website snapshots, researchers and developers can gain insights into the overall design aesthetic, User Interface (UI) elements, and User Experience (UX) considerations of e-commerce platforms. They can observe how products are showcased, how information is organized, and how users interact with different elements on the page.

Website snapshots can be utilized for various purposes, including:

- **Layout Analysis**: Researchers can analyze website snapshots to understand the overall layout structure of e-commerce websites. This includes identifying the placement of key elements such as product listings, featured items, promotional banners, and navigation bars.

- **Visual Element Detection**: Website snapshots enable the detection and analysis of visual elements present on the web page. This may include identifying product images, logos, text overlays, call-to-action buttons, and other graphical elements.

- **User Behavior Studies**: By studying website snapshots, researchers can gain insights into user behavior patterns, such as the areas of the page that attract the most attention, the path users follow while navigating the site, and the effectiveness of different design elements in influencing user interactions.

- **Competitive Analysis**: Website snapshots can be used for comparative analysis, allowing researchers to compare the design and layout of different e-commerce websites within the same industry or niche. This helps in identifying trends, best practices, and areas for improvement.

Overall, website snapshots serve as a valuable tool for analyzing the visual aspects and layout structures of e-commerce websites, providing insights that can inform design decisions, user experience optimizations, and competitive strategies.

## 3.2. OBJECT DETECTION (YOLO):

Object Detection, particularly with the You Only Look Once (YOLO) algorithm, is a cutting-edge technology that revolutionizes the way computers identify and locate objects within images or video frames. YOLO is renowned for its speed and accuracy, making it ideal for applications where real-time processing is crucial, such as autonomous vehicles, surveillance systems, and, in this case, product recognition in e-commerce.

The YOLO algorithm employs a single neural network to simultaneously predict bounding boxes and class probabilities for multiple objects within an image. Unlike traditional object detection algorithms that perform region proposals and classification separately, YOLO divides the input image into a grid and predicts bounding boxes and probabilities directly from the grid cells. This unique approach enables YOLO to achieve remarkable speed without compromising accuracy.

In the context of product recognition in e-commerce, YOLO can swiftly analyze product images and precisely locate various elements such as products, logos, text, or other relevant features. This capability is invaluable for tasks like inventory management, visual search, or content moderation on e-commerce platforms. By efficiently identifying and locating objects, YOLO streamlines the process, enhances user experience, and enables businesses to make data-driven decisions.

Furthermore, YOLO is highly versatile and can be trained on custom datasets to recognize specific objects tailored to the needs of a particular application. This adaptability makes it suitable for various industries and use cases beyond e-commerce, including robotics, medical imaging, and security.

Overall, YOLO object detection technology, with its swift and accurate identification and localization capabilities, empowers businesses to efficiently analyze product-related elements in images, ultimately enhancing processes and driving innovation in e-commerce and beyond.

Fig 3.1: YOLO Architecture

## 3.3. TEXT EXTRACTION (TESSERACT OCR):

Text Extraction using Tesseract OCR involves utilizing the Tesseract Optical Character Recognition (OCR) engine, a powerful open-source tool, to extract textual information from various sources, particularly images. Tesseract OCR works by analyzing the pixels of an image and identifying patterns that resemble characters. It then converts these patterns into machine-readable text, effectively extracting the textual content embedded within the image.

The process typically begins with preprocessing the image to enhance its clarity and remove any noise or distortions that may hinder accurate character recognition. Tesseract OCR then segments the image into individual characters and analyzes each character to determine its corresponding text. This involves recognizing the shapes and structures of letters, numbers, and symbols within the image.

One of the key advantages of Tesseract OCR is its ability to handle various fonts, languages, and writing styles, making it highly versatile for text extraction tasks across different types of images. It can recognize text in multiple languages and even handle complex scripts and languages with non-Latin characters.

Furthermore, Tesseract OCR is customizable and can be fine-tuned to improve accuracy and performance for specific applications or use cases. This customization may involve training the OCR engine with additional data or adjusting parameters to optimize its performance for particular image types or languages.

Overall, Tesseract OCR provides a robust and efficient solution for extracting textual information from images, enabling automation and digitization of text-heavy documents, forms, labels, and other visual content. Its versatility, accuracy, and open-source nature make it a popular choice for various text extraction applications in industries such as document management, data analysis, and automation.

Fig 3.2: OCR model Architecture

## 3.4. DATA PROCESSING:

Data processing is a critical step in the data extraction workflow, involving the transformation and organization of raw extracted data to ensure its cleanliness, remove noise, and structure it in a way that facilitates subsequent analysis.

**3.4.1. Cleaning Data:** The first aspect of data processing involves cleaning the extracted data to remove any inconsistencies, errors, or irrelevant information. This may include removing duplicates, correcting typos, standardizing formats, and handling missing values. Cleaning the data ensures that it is accurate and reliable for further analysis.

**3.4.2. Filtering Noise:** During extraction, irrelevant or noisy data may be captured along with the desired information. Data processing involves filtering out this noise to focus only on the relevant data. This could involve setting thresholds, applying statistical techniques, or using machine learning algorithms to identify and remove noise from the dataset.

**3.4.3. Structuring Information:** Once the data is cleaned and noise is filtered out, it needs to be structured in a way that makes it suitable for analysis. This involves organizing the data into a structured format such as tables, graphs, or hierarchical structures, depending on the nature of the data and the analysis requirements. Structuring the information makes it easier to interpret and analyze, enabling insights to be drawn more effectively.

**3.4.4. Normalization and Standardization:** In some cases, data processing also involves normalizing or standardizing the data to ensure consistency and comparability. This could include converting units of measurement, scaling numerical values, or standardizing categorical variables. Normalization and standardization make the data more uniform and facilitate meaningful comparisons across different datasets or time periods.

**3.4.5. Data Integration:** Data processing may also involve integrating data from multiple sources to create a unified dataset for analysis. This could include merging datasets, resolving

inconsistencies, and creating relationships between different pieces of information. Data integration ensures that all relevant data is considered in the analysis, providing a more comprehensive understanding of the underlying phenomena.

Overall, data processing is essential for preparing extracted data for meaningful analysis. It involves cleaning, filtering, structuring, and integrating data to ensure its quality, reliability, and usability for subsequent analytical tasks.



## Steps for data preprocessing

1 Data profiling
3 Data reduction
5 Data enrichment

2 Data cleansing
4 Data transformation
6 Data validation

Fig 3.3: Steps of Data preprocessing

## 3.5. BERT TRANSFORMATION MODEL

The BERT (Bidirectional Encoder Representations from Transformers) transformation model is a state-of-the-art Natural Language Processing (NLP) technique developed by Google. It utilizes a transformer architecture to pre-train deep bidirectional representations of text, enabling it to capture contextual information and semantic relationships within language data effectively. By leveraging the BERT transformation model, businesses can enhance their understanding of textual information related to products. BERT can analyze product descriptions, reviews, and other textual content associated with products, establishing semantic relationships between words, phrases, and sentences. This contextual understanding allows for more nuanced interpretation of text, capturing subtle nuances and connotations that traditional keyword-based approaches might miss.

The key advantage of using the BERT transformation model is its ability to comprehend the context in which words and phrases appear, rather than treating them in isolation. This contextual understanding enables more accurate and meaningful analysis of textual data, leading to improved insights and decision-making.

In the context of e-commerce or product-related applications, leveraging the BERT transformation model can facilitate tasks such as product categorization, sentiment analysis, and recommendation systems. By better understanding the textual information associated with

products, businesses can tailor their strategies more effectively, leading to enhanced customer experiences and improved business outcomes.

Example values for proposed method…



Fig 3.4. Web snapshots of Amazon website



Fig 3.5. Data objects detected by YOLO

```
File Name:C:\git\yolo\demo\input\uc2-wl-2.JPG
Title:The Power of Your Subconscious Mind
Title:RRB 2020 Maths (General and Advance) Chapter-wise & Type-wise Solved Papers
Price:1192
Author:by Youth Competition Times
Price:199
Author:by Joseph Murphy
```

Fig 3.6. Extracted Data from detected objects

## 3.6. TRANSFER LEARNING

Transfer learning is a technique in machine learning and deep learning where knowledge gained from training one model is transferred or applied to another related task or domain. In the context of data extraction from online sources, transfer learning can be employed in deep learning models to adapt to changes in website structures.

When a deep learning model is trained on a specific task, such as image classification or natural language processing, it learns to extract features from the input data that are relevant to that task. Transfer learning leverages this learned knowledge by fine-tuning the pre-trained model on a new, related task. In the case of data extraction from websites, the pre-trained model could have been trained on a similar task, such as web page parsing or information extraction.

By applying transfer learning, the deep learning model can adapt to changes in website structures without requiring manual intervention or re-training from scratch. This is particularly advantageous in scenarios where websites frequently update their layouts or structures, as the model can continue to efficiently extract relevant data even as these changes occur.

Transfer learning enables sustained efficiency in data extraction by allowing the model to leverage the knowledge gained from previous tasks, thus reducing the need for extensive re-training or re-engineering efforts. Additionally, it can improve the model's robustness to variations in website layouts and structures, enhancing its generalization capability across different domains or datasets.

Overall, transfer learning offers a powerful approach to adapting deep learning models for data extraction tasks from online sources, enabling more efficient and scalable solutions in the face of dynamic website environments.

## 3.7. FEATURE ENGINEERING

Feature engineering is a crucial aspect of data preprocessing in machine learning and data analysis. It involves the process of creating new features or transforming existing ones to enhance the dataset's predictive power and improve the performance of machine learning models. By crafting new features, feature engineering aims to capture meaningful patterns and relationships within the data, thereby providing more relevant and valuable information for analysis and prediction tasks. The process of feature engineering typically involves several steps:

**3.7.1. Understanding the Data:** Before crafting new features, it's essential to thoroughly understand the dataset, including its structure, variables, and the problem domain. This understanding helps identify potential opportunities for feature engineering and informs the creation of relevant features.

**3.7.2. Feature Selection:** Feature selection involves choosing the most relevant features from the dataset to include in the analysis. This step helps reduce dimensionality and computational complexity while retaining the most informative attributes for modeling. Feature selection can be performed before or after feature engineering, depending on the specific context and goals of the analysis.

**3.7.3. Creating New Features:** This is the core of feature engineering, where new features are crafted based on domain knowledge, intuition, or statistical techniques. New features may be derived from existing variables through mathematical transformations, combinations of variables, or domain-specific rules. For example, in a dataset containing information about housing prices, new features such as the ratio of bedrooms to bathrooms or the age of the property since its last renovation could be created to capture additional information relevant to predicting housing prices.

**3.7.4. Encoding Categorical Variables:** Categorical variables, such as gender or product categories, often require encoding into numerical values before they can be used in machine learning models. Feature engineering may involve converting categorical variables into numerical representations using techniques like one-hot encoding, label encoding, or target encoding.

**3.7.5. Handling Missing Values:** Missing values in the dataset can significantly impact model performance. Feature engineering may involve imputing missing values using techniques such as mean or median imputation, predictive modeling-based imputation, or encoding missingness as a separate feature.

**3.7.6. Scaling and Normalization:** Scaling and normalization techniques may be applied to numerical features to ensure that they are on a similar scale and have comparable ranges. This step helps prevent features with larger magnitudes from dominating the model training process and ensures that the model learns from each feature equally.

Overall, feature engineering plays a crucial role in extracting valuable insights from data and

building robust machine learning models. By creating new features that capture relevant information and preprocessing existing features appropriately, feature engineering enhances the dataset's quality and contributes to more accurate and effective data analysis and prediction.

## 3.8. DATA ANALYSIS

Data analysis is a crucial component of the data processing pipeline, involving the examination and interpretation of processed data to uncover patterns, trends, and valuable insights. This process aims to extract actionable intelligence that can inform decision-making and drive meaningful outcomes.

There are several key steps involved in data analysis:

**3.8.1. Exploratory Data Analysis (EDA):** This initial phase involves exploring the data to understand its structure, distribution, and characteristics. Techniques such as summary statistics, data visualization, and correlation analysis are commonly used to gain insights into the dataset.

**3.8.2. Pattern Recognition:** Once the data is understood, analysts employ various statistical and machine learning techniques to identify patterns and relationships within the dataset. This may involve clustering similar data points, detecting trends over time, or uncovering associations between variables.

**3.8.3. Insight Generation:** Data analysis aims to generate actionable insights that can drive decision-making. This involves interpreting the patterns and trends identified in the data to extract meaningful information relevant to the problem at hand. Insights may include identifying customer preferences, market trends, or areas for process optimization.

**3.8.4. Validation and Interpretation:** It's essential to validate the findings of the analysis to ensure their accuracy and reliability. This may involve conducting hypothesis testing, cross-validation, or sensitivity analysis to assess the robustness of the results. Interpretation of the findings involves translating the analytical results into actionable recommendations that stakeholders can understand and act upon.

**3.8.5. Visualization and Reporting:** Communicating the results of the analysis effectively is crucial for facilitating decision-making. Data visualization techniques such as charts, graphs, and dashboards are used to present the findings in a clear and understandable manner. Additionally, comprehensive reports are often generated to document the analysis process, results, and recommendations for further action.

Overall, data analysis is a multifaceted process that involves extracting valuable insights from processed data to drive informed decision-making and achieve desired outcomes. By leveraging advanced analytical techniques and tools, organizations can unlock the full potential of their data and gain a competitive edge in today's data-driven landscape.



Fig 3.7: E-Commerce data analysis

## 3.9. REPORT GENERATION

Report generation is a crucial aspect of data analysis, involving the compilation and presentation of key insights, trends, and findings extracted from the analyzed data. This process transforms raw data into actionable information, enabling stakeholders to make informed decisions.

The first step in report generation is to define the scope and objectives of the report. This involves identifying the target audience, determining the key metrics to be analyzed, and outlining the desired format and structure of the report.

Once the scope is defined, data analysis techniques are applied to extract relevant insights from the dataset. This may involve statistical analysis, data visualization, trend analysis, or machine learning algorithms, depending on the nature of the data and the specific objectives of the report.

After analyzing the data, the next step is to organize and structure the findings into a coherent narrative. This includes summarizing key insights, identifying trends and patterns, and highlighting any significant findings or anomalies.

The report should be structured in a clear and concise manner, with sections devoted to each key aspect of the analysis. Visual aids such as charts, graphs, and tables can be used to enhance the presentation of the data and make complex information more digestible.

In addition to presenting the findings, the report should also provide context and interpretation to help stakeholders understand the implications of the data. This may involve comparing the results to historical data, benchmarking against industry standards, or providing recommendations for future actions.

# 4. <u>IMPLEMENTATION</u>

## 4.1. Functionalities:

The system offers a comprehensive suite of functionalities aimed at revolutionizing the process of extracting product data from e-commerce websites. Users begin by selecting images containing the products they wish to analyze, initiating a seamless journey through automated data extraction. Leveraging YOLO (You Only Look Once) object detection, the system swiftly identifies and isolates individual products within the images. This process encompasses the precise detection of key attributes such as product images, prices, discounts, and offers associated with each item. Following object detection, the system employs sophisticated text extraction techniques to refine and format the extracted product information. This step ensures the removal of irrelevant details and the presentation of data in a clear and consistent manner. The culmination of these functionalities' manifests in the generation of structured datasets, stored in Excel sheets, primed for further analysis and decision-making. By streamlining the data extraction process, the system empowers users with actionable insights, facilitating efficient analysis and informed decision-making in the dynamic landscape of e-commerce.

### 4.1.1    Extracting Screenshots from the Browser:

It involves utilizing programming libraries or browser automation tools to capture images of web pages displayed within a browser window. This process typically entails the following steps:

Browser Automation

The system may utilize browser automation frameworks such as Selenium WebDriver or Puppeteer to control web browsers programmatically. These frameworks allow the system to navigate to specific URLs, interact with page elements, and captures screenshots.

Capturing Screenshots:

Once the browser is directed to the desired web page, the system triggers a command to capture a screenshot of the entire page or specific sections of it. This can be achieved using built-in browser functions or by invoking screenshot commands provided by the automation framework.

Image Storage:

The captured screenshots are then stored in memory or saved to a designated location on the system's file system. This ensures that the images are readily accessible for subsequent processing and analysis.

Metadata Collection:

Alongside capturing screenshots, the system may also collect metadata such as the URL of the web page, timestamp of the screenshot, and any relevant session information. This metadata provides context for the captured images and aids in organizing and analysing the data.

Error Handling:

The system incorporates error-handling mechanisms to address potential issues during the screenshot capture process, such as network errors, page load failures, or unexpected browser behaviour. This ensures the robustness and reliability of the screenshot extraction functionality.

The "Extracting Screenshots from the Browser" functionality involves programmatically controlling a web browser to capture visual snapshots of online shopping sites, enabling the acquisition of raw visual data for subsequent analysis and processing within the system.

### 4.1.2 Taking Product from the Screenshots:

Image Capture:

Initially, the system captures screenshots of web pages displayed in a browser. This typically involves using libraries or APIs provided by programming languages like Python to programmatically take screenshots of the browser window.

Image Processing:

Once the screenshots are captured, the system applies image processing techniques to analyze the images and locate regions containing product information. This can include methods such as edge detection, color segmentation, and object recognition to identify and isolate areas of interest within the screenshot.

Object Detection:

 After isolating regions of interest, the system employs object detection algorithms to   identify and extract individual product images from the screenshot. This  may involve using pre-trained deep learning models such as Convolutional Neural Networks (CNNs) to recognize common product shapes and features.

Feature Extraction:

 Once the product images are identified, the system extracts relevant features or descriptors from each image. This can include characteristics such as color histograms, texture patterns, or shape descriptors, which are used to represent the visual content of the products.

Product Classification:

Finally, the system classifies each extracted product image based  on its visual features. This classification process may involve comparing the extracted features to a database of known product categories or training a machine learning model to classify products into predefined categories.

### 4.1.3. Identifying product details in each product:

Feature Extraction:

After the products are extracted from the screenshots, the system uses feature extraction techniques to identify key attributes within each product image. These    attributes may include text, shapes, colors, and other visual elements that signify important product details.

Text Recognition:

For textual attributes such as product name, description, and brand, Optical Character Recognition (OCR) algorithms are employed to convert the text within the product images into machine-readable format.

Image Processing:

For non-textual attributes such as product images and availability indicators, image processing algorithms are applied to analyze visual features and extract relevant information. This may involve techniques such as object detection, image segmentation and pattern recognition.

Data Parsing and Classification:

Once the attributes are extracted, the system parses the data to identify and categorize specific details such as product name, price, description, brand, and availability. Classification algorithms are used to assign each attribute to its corresponding data field. Validation and Quality Assurance:

Finally, the extracted product details are validated and subjected to quality assurance checks to ensure accuracy and completeness. This may involve cross-referencing external databases or comparing with known product information to verify the correctness of the extracted details.


### 4.1.4. Extracting text for Each Product and Formatting:

Text Extraction:

Utilizing techniques such as Optical Character Recognition (OCR) or text parsing algorithms, the system extracts text data from the product details displayed on the webpage. This includes information like product descriptions, specifications, and other textual content associated with each product.

Text Formatting:

Once the text data is extracted, it undergoes formatting to ensure consistency and readability. This may involve tasks such as removing unnecessary characters, standardizing text formatting (e.g., font size, style), and organizing the text into a structured layout.

Identification of Data Fields:

The system identifies the specific data fields to which each extracted text belongs. This involves parsing the text to determine which pieces of information correspond                                      to attributes such as product name, price, description, brand, availability, and any other relevant details.

Mapping Text to Data Fields:

Each extracted text snippet is mapped to its corresponding data field based on predefined rules or patterns. For example, text mentioning the product's name may be assigned to the "product name" field, while text containing pricing information may be assigned to the "price" field.

Data Structuring:

Finally, the extracted text data, organized by data fields, is structured into a standardized format suitable for further processing or storage. This ensures that the extracted product

details are consistent and can be easily accessed and utilized for subsequent analysis or presentation.

Overall, this functionality enables the system to effectively extract text data from product details on webpages, format it into a structured layout, and accurately assign it to relevant data fields, facilitating the systematic organization and utilization of product information for various purposes.

## 4.2. ATTRIBUTES:

Utilization of YOLO object detection:
 The system harnesses the power of YOLO (You Only Look Once) object detection, renowned for its accuracy and efficiency in identifying objects within images.

Precise extraction of product attributes:
 With YOLO object detection, the system accurately extracts key attributes such as product images, prices, discounts, and offers, ensuring comprehensive data capture is essential for analysis.

Advanced text extraction techniques:
The system employs sophisticated text extraction techniques to refine the extracted product information, enhancing clarity and consistency in the data.

Robust design for versatility:
Designed to accommodate a diverse range of product images from various e-commerce websites, the system demonstrates robustness and versatility across different platforms and scenarios.

Adaptability for different use cases:
The synergy of these attributes culminates in a highly adaptable and effective system capable of extracting and organizing product data with unparalleled precision and efficiency.

## 4.3. DATASET:

Central to the system's functionality is the meticulously curated dataset used for training and testing purposes. Comprising two primary components—image data and labeled annotations—the dataset forms the foundation for the system's performance. The image data encompasses a diverse collection of approximately 500 images sourced from various e-commerce websites. These images span a wide spectrum of products and scenarios, providing a rich and nuanced dataset for model training. Each image undergoes meticulous manual labeling to identify and annotate individual products, specifying attributes such as price, image, discount, and offers. Additionally, a subset of approximately 450 images per product category is selected and labeled with detailed annotations, further enriching the dataset and facilitating rigorous evaluation. This comprehensive dataset serves as the linchpin of the system, enabling the robust training and fine-tuning of the object detection model for accurate product extraction.



Fig 4.1: CSV format of images and related data

## 4.4. Methodology:

The system's implementation is underpinned by a methodical and meticulous approach aimed at optimizing efficiency and reliability in extracting product data from e-commerce platforms. It initiates with a thorough data collection phase, where an extensive and diverse dataset of e-commerce product images undergoes careful curation and annotation. This dataset serves as the groundwork for training a cutting-edge YOLO (You Only Look Once) object detection model, renowned for its swift and precise identification of objects within images. Through iterative training on the annotated dataset, the model hones its abilities, progressively improving its capacity to identify and extract product information from various image types. Once satisfactorily trained, the YOLO model is deployed to process user-selected images from e-commerce sites. Utilizing its sophisticated algorithms and neural network architecture, the model swiftly analyzes these images, capturing essential product attributes such as price, image, discount, and promotions. This phase

is crucial for ensuring the extracted data's comprehensiveness, accuracy, and timeliness, providing businesses with invaluable insights into market trends and product offerings. After meticulously curating and annotating a diverse set of e-commerce product images, the system trains a YOLO (You Only Look Once) object detection model. This model is optimized to swiftly and accurately recognize product attributes like price, image, discount, and offers. Through rigorous training on the labeled dataset, the YOLO model refines its ability to extract pertinent object details for each product depicted in the images. Once trained, it efficiently analyzes product images, providing optimized and accurate object details essential for e-commerce analysis and decision-making. Subsequently, the system employs advanced text extraction techniques to further enhance the extracted product information. By parsing through textual elements within images, such as product descriptions and specifications, relevant data is extracted and integrated into the dataset, enhancing its clarity and usefulness for subsequent analysis and decision-making. Finally, the refined data is organized and stored in Excel sheets, structured in a manner conducive to further analysis and interpretation. This centralized data repository serves as a valuable asset for businesses, enabling strategic planning, targeted marketing campaigns, and pricing optimization. In essence, this systematic approach underscores the system's commitment to precision, efficiency, and reliability in automating product data extraction, providing businesses with actionable insights for informed decision-making and competitive advantage. Moreover, the system places a high priority on data security and privacy, implementing robust encryption protocols and access controls to safeguard sensitive information. Continuous monitoring and auditing mechanisms ensure compliance with regulatory standards, instilling confidence in users regarding the integrity and confidentiality of their data. Additionally, the system is designed for seamless integration with existing e-commerce platforms and systems, minimizing disruption to established workflows and facilitating efficient utilization of its capabilities.
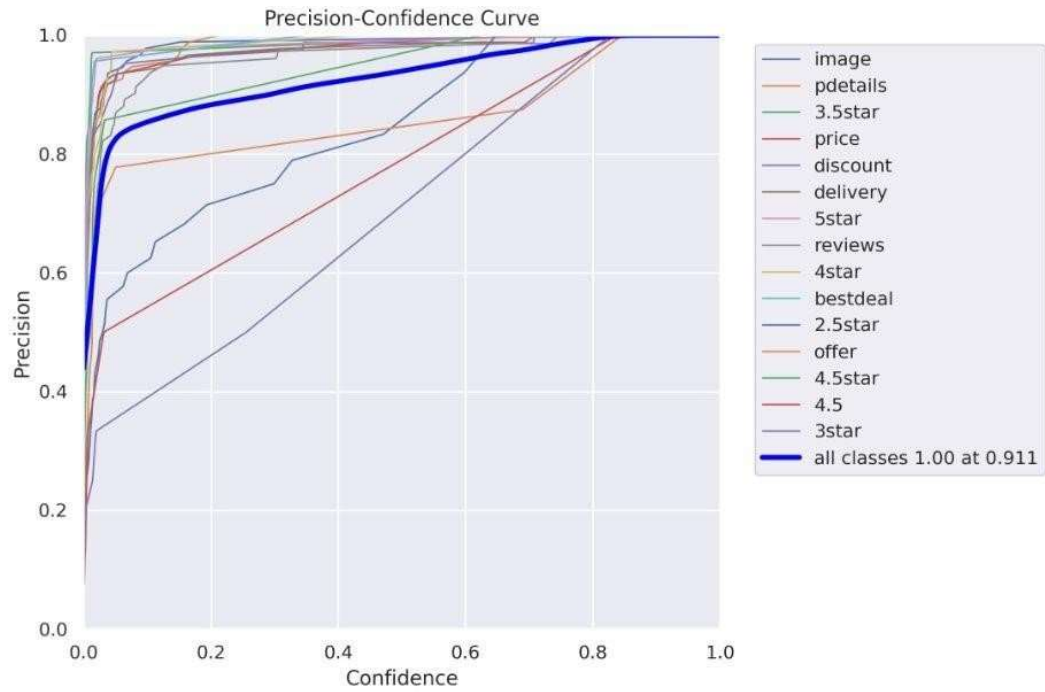
## 4.5. EXPERIMENTAL SCREENSHOTS:

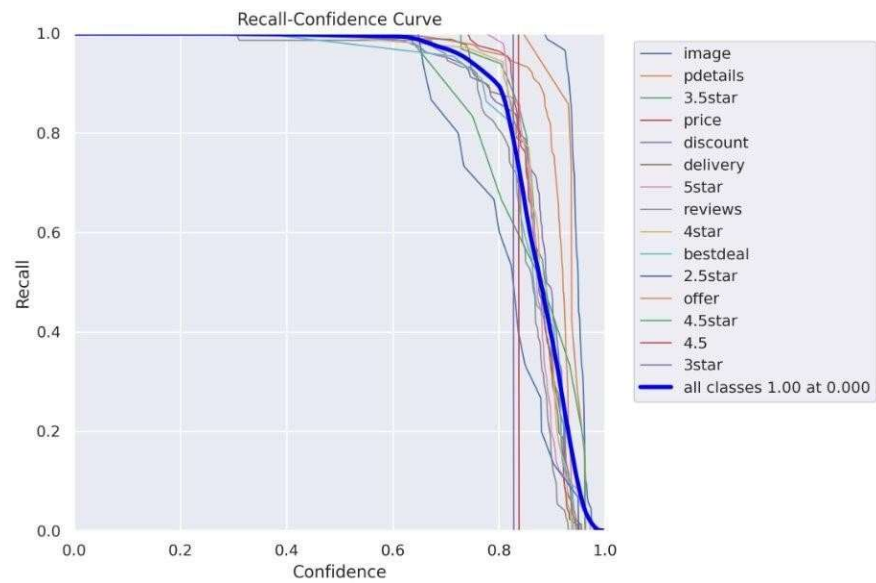

Fig 4.2. Snapshot of Precision graph
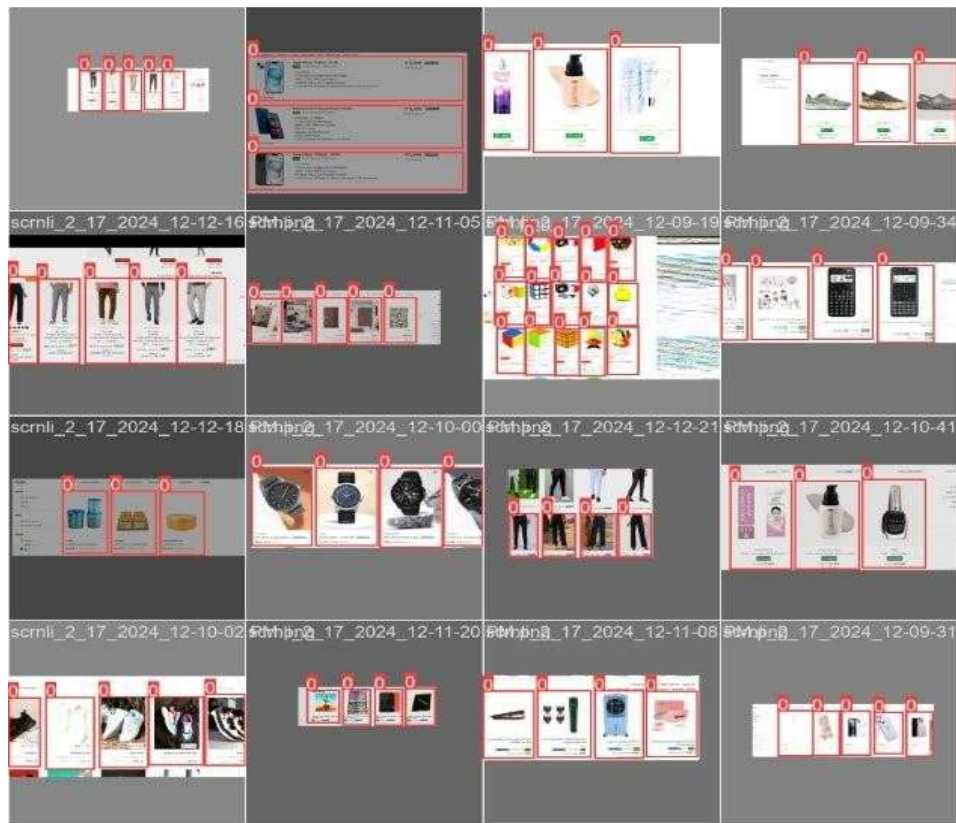


Fig 4.3. Snapshot of Recall graph

Fig 4.4. Identifying the objects

# 5. <u>EXPERIMENTAL RESULTS/ OBSERVATION</u>

## 5.1. Experimental Setup:

### 5.1.1. Setup Jupyter Notebook:

To install and set up Jupyter Notebook, you can follow these steps:

1. Install Python: First, you need to have Python installed on your system. You can download and install Python from the official Python website: https://www.python.org/. Make sure to check the option to add Python to your system PATH during installation.

2. Install Jupyter Notebook: Once Python is installed, you can install Jupyter Notebook using pip, which is the Python package manager. Open a terminal or command prompt and run the following command:  **pip install jupyter**

3. Launch Jupyter Notebook: After the installation is complete, you can launch Jupyter Notebook by running the following command in your terminal or command prompt: **jupyter notebook**

4. Accessing Jupyter Notebook: Once you run the command, your default web browser should open, and you'll be directed to the Jupyter Notebook dashboard. If it doesn't open automatically, you can manually open your web browser and go to http://localhost:8888/. Here, you'll see a file browser where you can navigate your filesystem and create or open Jupyter Notebook files.

5. Creating a New Notebook: To create a new notebook, click on the "New" button in the top right corner and select "Python 3" (or any other available kernel you want to use).

6. Using Jupyter Notebook: You can now start using Jupyter Notebook. Each notebook consists of cells where you can write and execute Python code, Markdown for documentation, and more. You can execute iby clicking on "Run" button or by pressing Shift+Enter.

7. Saving and Closing:  Make sure to save your work regularly by clicking the "Save" button or using the keyboard shortcut Ctrl+S. To close Jupyter Notebook, you can simply close the browser tab or stop the Jupyter Notebook server by pressing Ctrl+C in the terminal or command prompt where it's running.



Fig 5.1: Jupyter Notebook

## 5.1.2. Setting Up Visual Studio Code

To installing and configuring Visual Studio Code (VS Code) for a smooth development experience:

1. Download and Install:

Head to the official VS Code download page: https://code.visualstudio.com/download

Choose the installer suitable for your operating system (Windows, Mac, or Linux).

Run the downloaded installer and follow the on-screen instructions.

2. Open VS Code:

Once installed, locate and launch VS Code from your applications list.

3. Explore the Interface (Optional):

Take some time to familiarize yourself with the layout. Key areas include:

Activity Bar: Manage projects, extensions, and the terminal.

Side Bar: Explore opened folders and files.

Editor: The primary workspace for writing

code.

Panel: Provides additional information like output or version control.

Menu Bar: Access various settings and actions.

4. Install Extensions (Optional):

VS Code is powerful with extensions. Explore the Extensions marketplace within VS

Code and install language-specific extensions for syntax highlighting, code completion, and

debugging support relevant to your project needs.

5. Open Your Project:

Navigate to your project folder using the File Explorer within VS Code (File > Open

Folder).

6. Start Coding!

Now it is ready to start coding! Use the built-in features like syntax highlighting, code

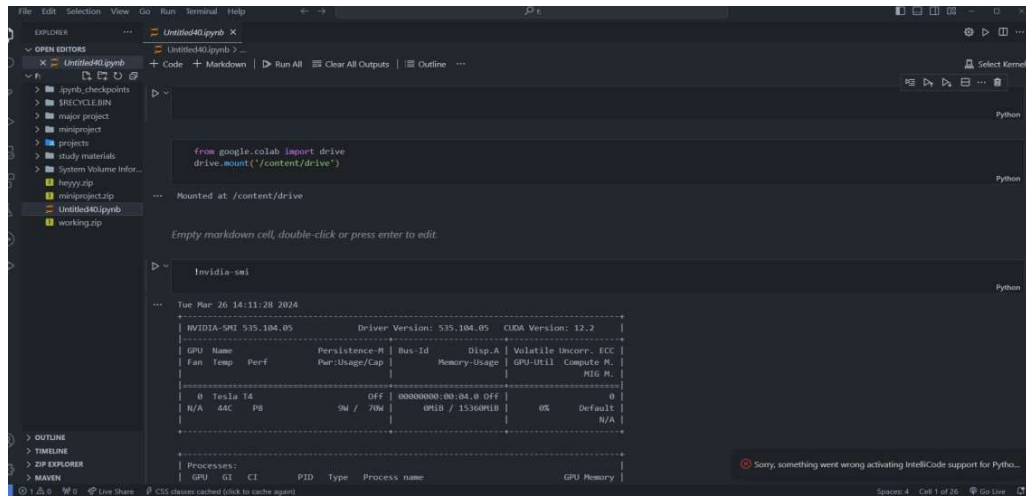completion, and debugging to write and test your code efficiently.



Fig 5.2: Visual Studio Code

## 5.2. PARAMETERS WITH FORMULAS

**Text Extraction with Tesseract OCR:**

Accuracy: The accuracy of Tesseract OCR can be measured using metrics such as Word Accuracy

or Character Accuracy.

- Word Accuracy:

  Word Accuracy = Number of Correctly Recognized Words Total Number of Words x

- Character Accuracy:

Character Accuracy = Number of Correctly Recognized Characters Total Number of Characters × 100%

Confusion Matrix: A confusion matrix can also be used to evaluate the performance of Tesseract OCR, providing insights into true positives, false positives, true negatives, and false negatives.

**OBJECT DETECTION:**

1.Intersection over Union (IoU):

IoU measures the overlap between predicted bounding boxes and ground truth bounding boxes. It is calculated as the ratio of the area of overlap to the area of union between the predicted and ground truth bounding boxes.

IoU = Area of Overlap / Area of Union

1. Precision and Recall:

Precision: Precision measures the proportion of true positive detections out of all positive predictions. It is calculated as the ratio of true positives to the sum of true positives and false positives.

Precision = True Positives / True Positives + False Positives

Recall (Sensitivity): Recall measures the proportion of true positive detections out of all ground truth objects. It is calculated as the ratio of true positives to the sum of true positives and false negatives.

Recall = True Positives / True Positives + False Negatives

2. Average Precision (AP):

AP summarizes the precision-recall curve into a single value. It is calculated as the area under the precision-recall curve.

3. Mean Average Precision (mAP):

MAP is the mean AP across all classes.it is calculated as the average of the AP values for all classes.

# 6. <u>DISCUSSION OF RESULTS</u>

The discussion of results encompasses an analysis of the YOLO model's performance, including its train and validation batch results, metrics, losses, and the generated confusion matrix. These outcomes provide valuable insights into the model's accuracy, precision, recall, and F1-score, shedding light on its ability to detect and classify objects effectively. Furthermore, the final analytics of product recognition offer a comprehensive overview of the model's performance, highlighting its strengths and areas for improvement. This discussion serves as a critical step in evaluating the YOLO model's efficacy and guiding future optimizations for enhanced object detection in real-world applications.
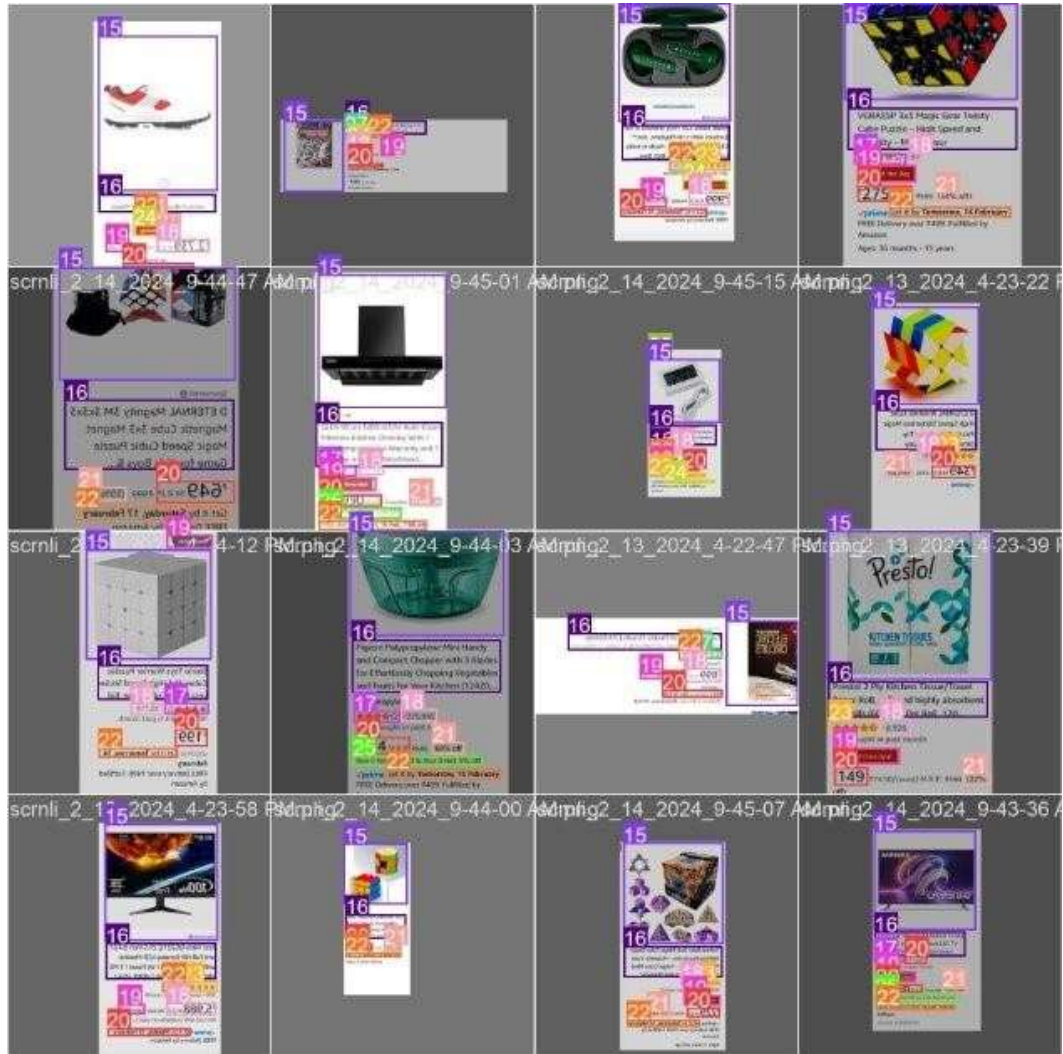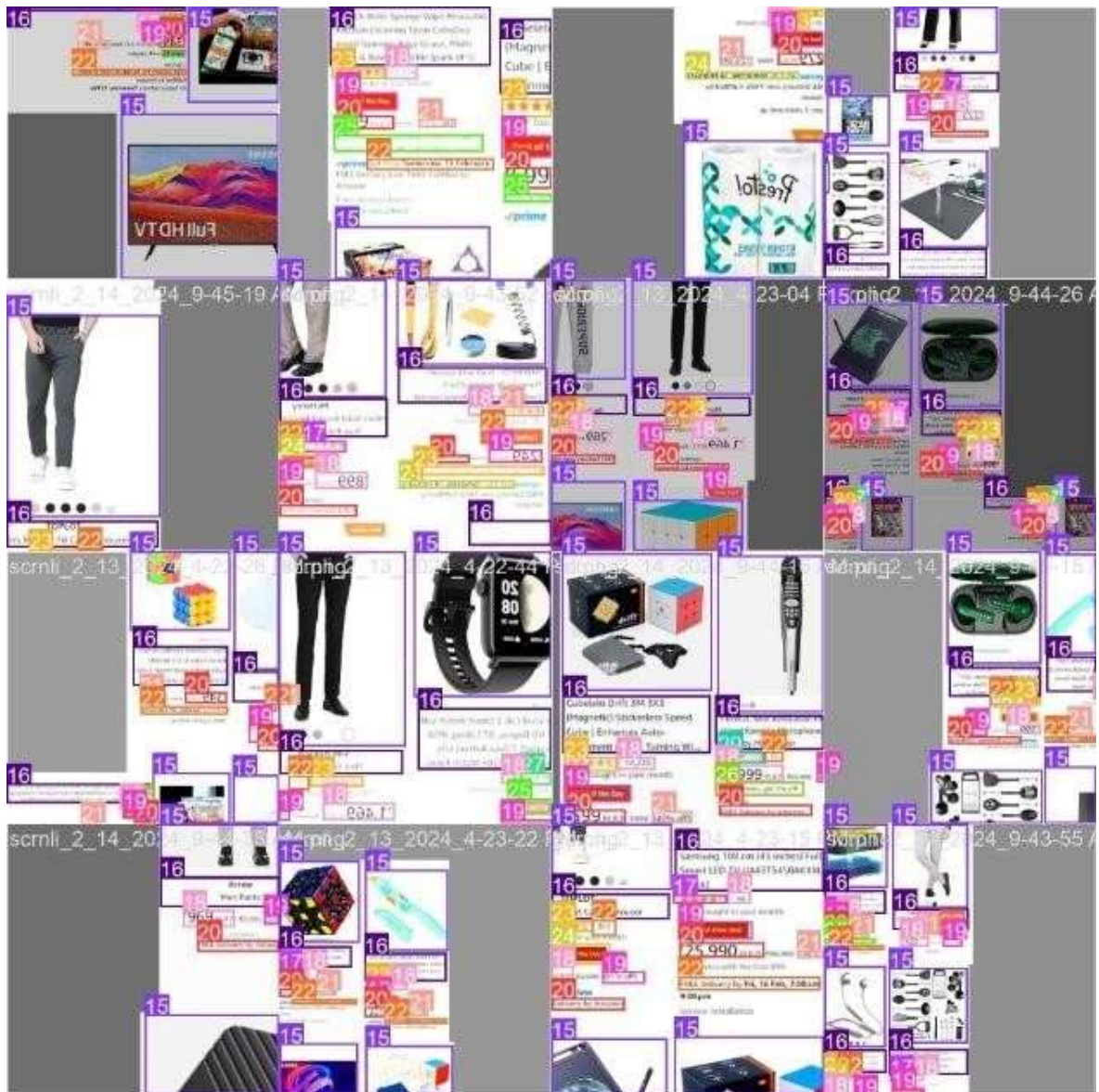


Fig 6.1. Train Data detection

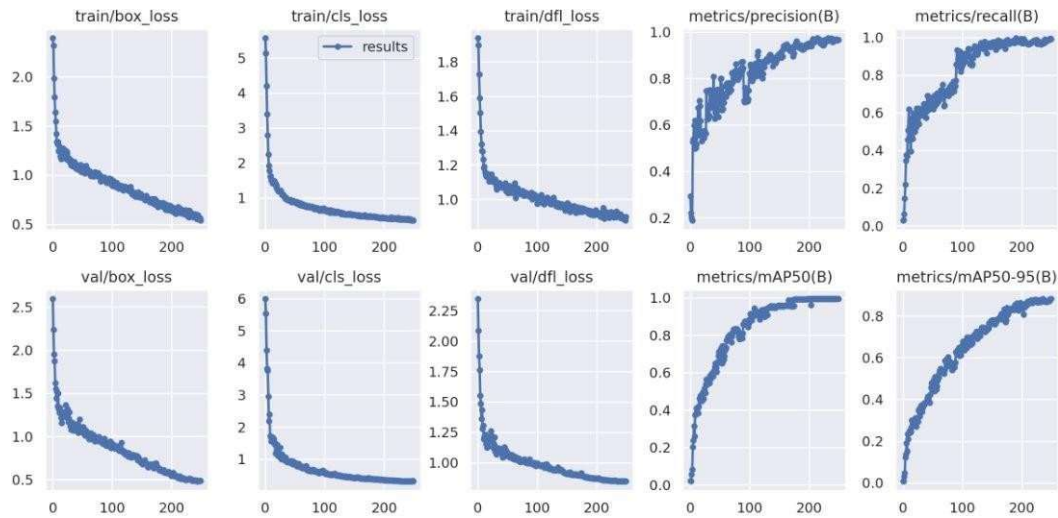Fig 6.2. Validation Data detection
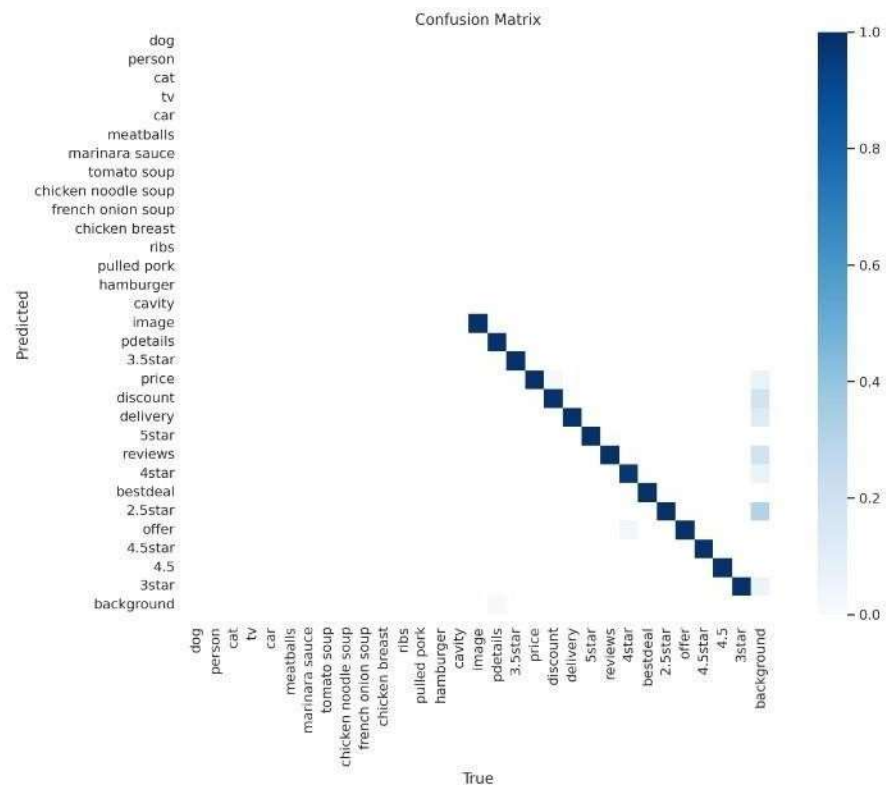
Fig 6.3. Losses and metrics results


Fig 6.4. Confusion matrix

# 7. SUMMARY, CONCLUSION AND RECOMMENDATION

This project developed an automated system for extracting structured product details, pricing, and reviews from e-commerce websites, boosting user trust and satisfaction. It streamlines inventory management, order fulfillment, and targeted marketing. To improve accuracy, especially with similar items, expanding training data and integrating computer vision techniques are proposed. The focus is on refining dynamic data extraction in e-commerce to enhance precision, engagement, and conversions, and to adapt to market changes for sustained growth and customer loyalty. Leveraging advanced algorithms and technologies like web scraping and API integrations ensures platforms deliver up-to-date, reliable information, fostering user trust and satisfaction while minimizing errors. Real-time data monitoring enables swift adaptation to market shifts, while personalized recommendations based on individual preferences enhance user experience, driving increased engagement and conversions. By prioritizing accuracy and user-centricity, e-commerce platforms can differentiate themselves and foster sustainable growth and customer loyalty. Our project aims to tackle the challenge of acquiring accurate product details from online shopping sites by integrating text comprehension and image recognition capabilities. This adaptive system enhances the online shopping experience by swiftly delivering reliable information, combining web scraping and API integration for data gathering with text comprehension and image recognition for processing. Incorporating adaptation strategies and reinforcement learning ensures flexibility and adaptability in dynamic online environments, marking a significant advancement in e-commerce data extraction efficiency and effectiveness, benefiting businesses and consumers alike.

# 8. <u>FUTURE ENHANCEMENT</u>

In future enhancements for the Dynamic E-commerce Data Extraction project, advanced machine learning and natural language processing algorithms can be integrated to enhance data extraction accuracy and efficiency. By employing techniques such as image recognition and sentiment analysis, the system can automatically extract product attributes from images and analyze customer reviews for sentiment and feedback. Additionally, the implementation of deep learning models can enable the system to adaptively learn from new data sources and continuously improve its extraction capabilities over time. Another area of enhancement involves expanding the scope of data analytics and insights provided by the system. Advanced data visualization techniques can be utilized to present analytics in a user-friendly and interactive manner, enabling businesses to gain deeper insights into consumer behavior, market trends, and competitor analysis. Moreover, predictive analytics algorithms can be implemented to forecast sales trends, identify emerging product categories, and optimize pricing strategies. Additionally, integration with external data sources such as social media platforms and industry reports can enrich the analytics provided by the system, offering a more comprehensive understanding of the e-commerce landscape. This includes analyzing social media mentions, influencer marketing effectiveness, and industry benchmarks to inform strategic decision-making. Overall, these future enhancements aim to further automate and optimize the e-commerce data extraction process while providing actionable insights and intelligence to businesses, ultimately improving their competitiveness and decision-making capabilities in the dynamic e-commerce market landscape.

# 9. REFERENCES

[1]. Sudhir Kumar Patnaik., "Intelligent and Adaptive Web Data Extraction System Using Convolutional and Long Short-Term Memory Deep Learning Networks," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., Dec 2021, v.4.

[2]. N. Islam, Z. Islam, and N. Noor, A survey on optical character recognition system, Journal of Information & Communication Technology-JICT, vol. 10, no. 2, pp. 1–4, 2016

[3]. J. Tao, H. B. Wang, X. Y. Zhang, X. Y. Li, and H. W. Yang, an object detection system based on YOLO in traffic scene, in Proc. of 2017 6th Int. Conf. Computer Science and Network Technology (ICCSNT), Dalian, China, 2017, pp. 315–319.

[4]. E. Uzun, A novel web scraping approach using the additional information obtained from web pages, IEEE Access, vol. 8, pp. 61726–61740, 2020

[5]. H. Rao and D. R. M. Sashikumar, A survey on automated web data extraction techniques for product specification from e-commerce web sites, International Journal of Advanced Research in Computer Science and Software Engineering, vol. 6, no. 8, pp. 310–316, 2016.

[6]. M. Salah, B. Al Okush, and M. Al Rifaee, A comparison of web data extraction techniques, in Proc. of 2019 IEEE Jordan Int. Joint Conf. Electrical Engineering and Information Technology (JEEIT), Amman, Jordan, 2019, pp. 785–789

[7]. S. L. Li, C. Chen, K. W. Luo, and B. Song, Review of deep web data extraction, in Proc. of 2019 IEEE Symp. Series on Computational Intelligence (SSCI), Xiamen, China, 2019, pp. 1068–1070.

[8]. C. J. Liu, Y. F. Tao, J. W. Liang, K. Li, and Y. H. Chen, Object detection based on YOLO network, in Proc. of 2018 IEEE 4th Information Technology and Mechatronics Engineering Conf. (ITOEC), Chongqing, China, 2018, pp. 799–803

[9]. S. Nagarajan and K. Perumal, A deep neural network for information extraction from web pages, in Proc. of 2017 IEEE Int. Conf. Power, Control, Signals and Instrumentation Engineering (ICPCSI), Chennai, India, 2017, pp. 918–922

[10]. T. Gogar, O. Hubacek, and J. Sedivy, Deep neural networks for web page information extraction, in Artificial Intelligence Applications and Innovations. IFIP Advances in Information and Communication Technology, vol. 475, L. Iliadis and I. Maglogiannis, eds. Thessaloniki, Greece: Springer, 2016, pp. 154–163

[11]. R. Girshick, J. Donahue, T. Darrell, and J. Malik, Rich feature hierarchies for accurate object detection and semantic segmentation, arXiv preprint arXiv: 1311.2524v5, 2014

[12]. D. Freitag, Information extraction from HTML: Application of a general machine learning approach, in Proc. of 15th National/Tenth Conf. Artificial Intelligence/Innovative Applications of Artificial Intelligence, Madison, WI, USA, 1998, pp. 517–523.

[13]. Y. Wang, A new concept using LSTM Neural Networks fordynamic system identification, in Proc. of 2017 American Control Conf. (ACC), Seattle, WA, USA, 2017, pp. 5324–5329.

[14]. E. Ferrara, P. De Meo, G. Fiumara, and R. Baumgartner, Web data extraction, applications and techniques: A survey Knowledge-Based Systems, vol. 70, pp. 301–323, 2014.

[15]. Y. H. Zhai and B. Liu, Web data extraction based on partial tree alignment, in Proc. 14th Int. Conf. World Wide Web, Chiba, Japan, 2005, pp. 76–85.

[16] S. Kumari and C. N. Babu, Real time analysis of social media data to understand people emotions towards national parties, in Proc. of 8th Int. Conf. Computing, Communication and Networking Technologies (ICCCNT), Delhi, India, 2017, pp. 1–6.

[17] D. G. Gregg and S. Walczak, Adaptive web information extraction, Communications of the ACM, vol. 49, no. 5, pp. 78–84, 2006.