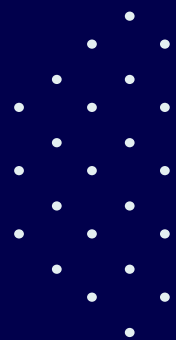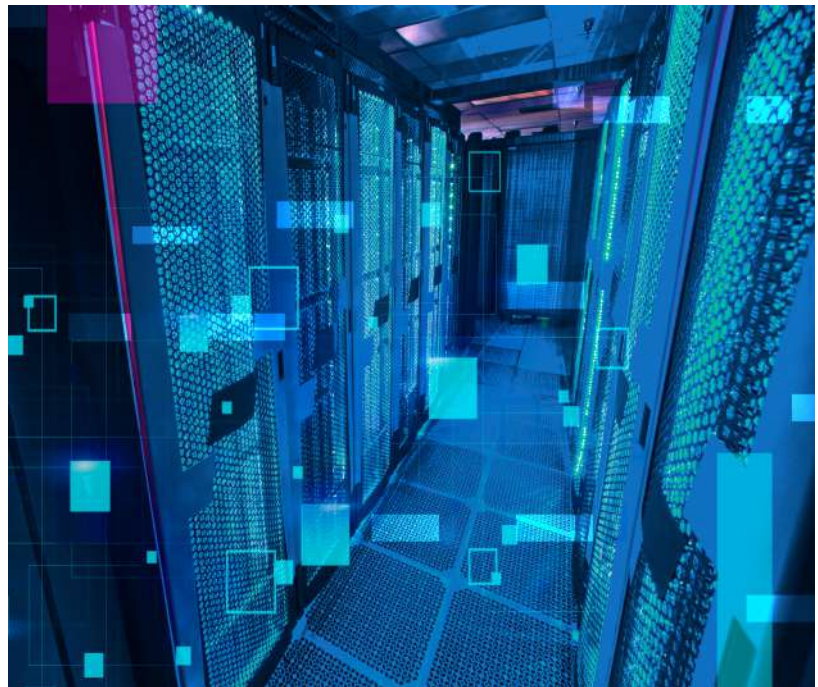Provectus

# Synthetic Invoice Dataset Generator

# Abstract

In this paper, we present the implementation of an invoice generator, with random filling of data and various template formats. This type of data can meet the needs of a machine learning system, to automate and improve the flow of complex document processing. Using a synthetic dataset solves the problem of sensitivity of the training data for such models. In addition, this invoice generator solves the problem of data annotation by automating the generation process. The quality of the generated dataset was checked using a custom model for information extraction. PDF or PNG formats, and output annotations to JSON or XML formats. The final dataset consists of invoices with varied ground truth and layout.

# Introduction

A synthetic dataset is a type of data that is generated by a program, and not collected from real life experience. Its main goal is to be flexible and rich enough to help conduct research with machine learning models.

In the process of building a machine learning model for data recognition and extraction from invoices, it is necessary to have a sufficiently large annotated dataset for the training process. Due to the sensitivity of information, such datasets are not publicly available, so we decided to build a document generator similar in format and content to examples of real-life invoices. This approach can be used not only for invoices, but also for other types of documents.

We propose the approach of building HTML invoice templates, along with further extraction of bounding box coordinates and labels during rendering. Our generator allows output documents to be converted to PDF or PNG formats, and output annotations to JSON or XML formats. The final dataset consists of invoices with varied ground truth and layout.

# Invoice Generation

After examining examples of invoices in real life, a typical structure for this type of document was formed, which is in three parts: the Head contains company and client information, invoice number and date information; the Central part contains tables with product information; and the Footer contains information about payment, total cost and taxes. Other information blocks are optional and can have different positions in the document layout.

To create an instance of an invoice, the information blocks of interest were first determined, for the extraction of values for which the model would be trained. For the field content of information blocks, random data were used for text fields and mathematical functions were used to calculate some numerical values. To generate the layout, random behavior of information blocks were organized in predetermined sections of the invoice structure.

There are several steps involved in creating invoices: creating HTML, extracting bounding box coordinates, extracting label data, and converting HTML files to PDF or PNG formats.

# Data Generation

To generate random text data we used the Python package Faker. Two prepared synthetic datasets were used for company names (Provectus Contract Generator) and product tables. Some of the fields were determined through functional dependencies, for example, the calculation of total cost and due date values.

Based on the analysis of real-life documents of this type, groups of possible names of target fields were collected to account for diversity in training. For example, "Invoice Date," "Issue Date," "Date," "Date of Issue," and "Issued on," can all represent the same target field in different documents.

# Layout Generation

The document layout was built with Jinja as a templating engine for Python, which allows for dynamic creation of HTML documents. Information blocks were defined for each part of the template (Head, Central, Footer). The rendering of all information blocks was randomized using specific css properties (background color, font size, border, etc.). Mixed render layouts were achieved in this way (Figure 1).
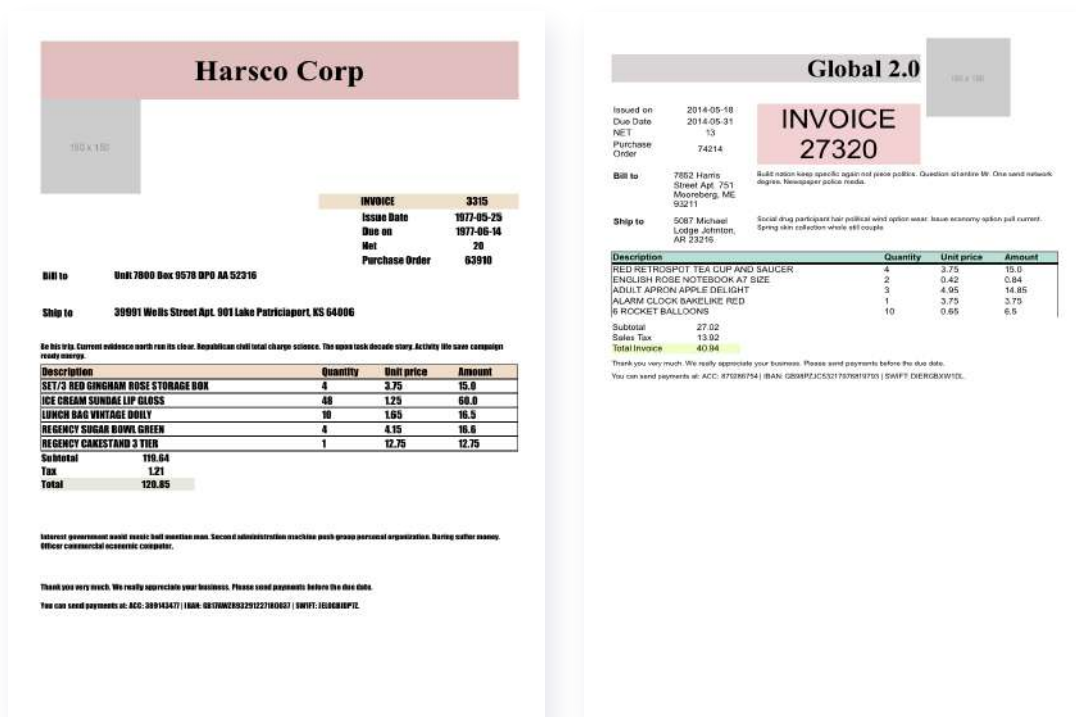
Figure 1. Examples of generated invoices.

# Annotation Generation

The annotation is included in a JSON or XML stream and contains the position and box size of the text zone, the text itself, and the annotation class that we want to learn.

# Invoice Analysis

In order to evaluate the performance of a machine learning model trained on data from our generator, we used a custom model based on a recently published article [1]. The key idea of the model is to extract information using knowledge about the types of target fields, to create candidates for extraction. The neural network architecture then studies the dense representation of each candidate based on neighboring words in the document. These learned representations are useful in solving the problem of extracting unseen document templates.

Initially, the image is run through the OCR engine (Textract or Tesseract). The text and the positions of the bounding boxes are extracted from the output. The text parser and predefined rules then extract the target field candidates and their neighboring words. Entities embedded using a word embedding table become input data for the neural network.

A generated synthetic dataset with 5000 images was used for training purposes. Evaluation of model performance was done using a dataset of real invoices with field result accuracy: invoice_id 0.68, invoice_date 0.72, total_amount 0.63, due_date 0.89, payment_terms 0.74, purchase_order 0.93, tax 0.78.

# Conclusion

This paper presents a method for creating a synthetic dataset represented as invoices. All documents have a unique render layout and design, and contain the variability of names of key fields, to approximate the documents to real examples of invoices.

As a result of this work, an artificial dataset for the purpose of preserving privacy was created that is able to create training data for machine learning algorithms. This work allowed us to create annotated sample invoices adequate for the learning and extraction of necessary information.

# References

Representation Learning for Information Extraction from Form-like Documents, Bodhisattwa Prasad Majumder, Navneet Potti, Sandeep Tata, James B. Wendt, Qi Zhao, Marc Najork, 2020

**Provectus**

# Synthetic Invoice Dataset Generator