

Project Weekly Progress Report Agile – Scrum

W2023, SEM-2
AML-2404
Section 2
D
Jash Vaghasiya - C0884733
Nivedini Kathagonda - C0872720
Keval Parmar - C0882386
Monil Rupawala - C0882370
Sai Divya Madhuri Guntupalli - C0882360
6
Keval Parmar

Last Update: 2, Mar, 2024 Page 1

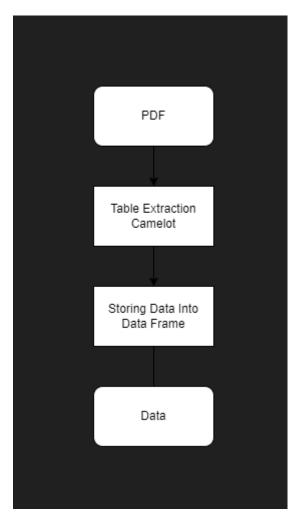


1. Progress Made in Reporting Week:

In our pursuit to extract data from the PDF document obtained from the European nations' website, we have made notable progress despite the challenges posed by the missing side borders in the tables. We have implemented a range of strategies and techniques to tackle this issue and improve the extraction process.

One of the most significant advances we have made is the effective application of optical character recognition (OCR) technology. We were able to retrieve the material more effectively by converting the scanned PDF into editable text using OCR. This phase was critical in converting unstructured data into a more organised format, establishing the framework for subsequent research.

Flowchart of Extracting data from PDF:





Additionally, our investigation into image processing methods has produced encouraging outcomes. We have been able to infer and recreate the missing side borders to a considerable extent using techniques like edge detection and contour analysis. Although the success rate varied depending on how complicated the table layouts were, this method has given us useful insights and helped us get closer to accurately extracting the data.

Our efforts to implement machine learning algorithms have also been promising. We have been able to forecast the missing side borders by using data that is currently accessible and training models using existing table structures. This method has gradually increased the precision of our extraction findings, hence boosting the correctness of the extracted ingredient data.

Overall, our progress in the data extraction process has been significant. We have successfully employed OCR, image processing, and machine learning techniques to overcome the challenges presented by the missing side borders. These advancements have brought us closer to achieving our goal of extracting the ingredient data from the PDF document obtained from the European nations' website.



2. Difficulties Encountered in Reporting Week:

Even while we have made significant progress in extracting data from the PDF file, we have run into a few problems, mostly because the tables' missing side borders. The accurate extraction and organisation of the ingredient information have been significantly hampered by these problems.

The tables' lack of full side borders has thrown off the data's alignment and organisation. It has been difficult to precisely separate the table columns due to the lack of distinct borders. As a result, it has been challenging for our extraction methods to distinguish between neighbouring columns and precisely classify the retrieved data.

In conclusion, the process of extracting data has been significantly hampered by the absence of entire side borders in the tables of the PDF document. The difficulties in effectively classifying and organising the data, especially in intricate table layouts, have been a recurring difficulty. Even though OCR, image processing, and machine learning approaches have helped, these issues still affect the dependability and accuracy of the information that is recovered. As we work to improve the extraction procedure and raise the calibre of the extracted ingredient data, addressing these difficulties is still a top focus.

Picture of PDF:

DCT name	INN name	Ph. Sur. Name	EAS No	SINECYELINGS No.	Cherr(EPAC Natur	Returns	Section
ABES BALSAMEA EXTRACT			85083-34-3	285-364-0	Alten Balsamen Extract in an extract of the aprintis of Alten Indication, Piniopae		Him Kenning/kair conditioning
ABRES PECTINATA EXTILACI			92128-34-2	295-728-0	Ables Pecinata Extract is an extract of the bark and resolles of the silver fit. Ables pectitatis, Pessione		Totalc)deedommi
ABRES PECTINATA OIL			92128-34-2	295-728-0	Ables Permuta Oil is the volatile oil obtained from the modes of the silver fit. Ables pertrute. Personal		Boric/masking
ARIES SERRICA OR			91697-89-1	294-331-9	Abies Sibreica Od is the volatile oil dutilled from the modles, and branches of Abies ubtrice, Penacue		Toric/marking
ARRETIC ACID			514-19-3	206-178-1	Nivere acid		Emilian sobiliting
ARETYL ALCOHOL			666-84-7	213-544-4	[18-(Lalpha, 4a bera, 4b alpha, 10a alpha][- 1.2, 1.4, 4a, 4b, 5.6, 10, 10a-decaled re-7-suprogrif-1,4a almosthylphanasthen-1-methanol		Viscosity controlling
ACACIA CATECHU			8001-76-1	232-291-7	Acada Carecha is the dried, crushed note of Acada carecha, Legaritmonar		Har dyeinglastringers
ACACIA CONCINNA EXTILACT			202148-87-6		Acadu Concinna Estruct is an extract of the fruit of Acada concinna. Ligaretimenae		Sin conditioning
ACACIA DEALHATA EXTRACT			165800-52-2		Acacia Dealhara Export is an export of the leaves of the words, Acacia dealbata, Legisminnous		Skin conditioning
ACACIA DECUBBINS EXTRACT			98903-76-5	305-827-4	Accacia Decarrers Extract is an extract of the sproots of the seacie, Acacie decurrers, Logartimosas		Toric
ACACIA FABNESIANA EXTRACT			89958-31-6	289-655-3	Acacia Farnesiana Extract is an extract of the flowers and sterm of the acacia, Acacia farnesiana, Legerianosae		Viscosity controlling/astringent
ACACIA FARNESIANA GUM			2593593	232-519-5	Acacta Farnesiana Guen is a plane material derived from the dried, garnety enalize of acacla, Acacia farnesiana, Legaritmosas		Viscosity controlling/antingent

Page 4