

## SCENARIO 1:

OCR **Dataset Creation** 

Data collection "OCR" Optical Character Data Analysis Recognition is a powerful technological process that

> The program behind this technology starts by analyzing the structure of the image file, it tries to separate the different structures that compose the image: tables,

texts, images, etc.

converts a paper document,

PDF or even a photo into an

Automated data extraction

editable digital document.

Then each line is converted into words and then into characters. Once the characters have been defined, the OCR begins to compare them with a group of images already processed by the OCR or predefined which allow it to propose hypotheses on the meaning of these characters.

The OCR uses these assumptions to analyze the curvatures of lines in words and words in characters. After reviewing all the assumptions, the program ultimately proposes content that could be consistent with the initial

Format error detection (optional)

> Format errors, in the best-case scenario, can break an automated data processing pipeline. In the worst case, they introduce logical errors in downstream analytical tasks that are difficult to debug.

Example: international different way to write dates.

In our case, the documents are received in different formats.

Regex method assessment for SDS

=> Focus on the Safety Data **Sheet** 

Regular Expression – better

known as Regex – is pattern

matching. They are used to

match well-defined patterns (= matching tools identify words and phrases with common patterns). It is not an Al solutions using Natural Language Processing (NLP) and Machine Learning (ML) that can conduct an intelligent analysis of unstructured text. Indeed, NLP understands the language as it allows computers to undertake linguistic analysis and in practical terms "read" a document just like a human. In our situation, we don't need to understand the text. We

need to report the content of

document which are always

following the same pattern.

Apply the Regex model on SDS

Apply the SDS regex model on all SDS documents.

Regex method assessment for TD

=> Focus on the Test Report

Regular Expression – better

known as Regex – is pattern matching. They are used to match well-defined patterns (= matching tools identify words and phrases with common patterns). It is not an Al solutions using Natural Language Processing (NLP) and Machine Learning (ML) that can conduct an intelligent analysis of unstructured text. Indeed, NLP understands the language as it allows computers to undertake linguistic analysis and in practical terms "read" a document just like a human. In our situation, we don't need to understand the text. We need to report the content of document which are always

following the same pattern.

Apply the Regex

model on TR

Apply the Regex model on all

Test Report (TR) documents.

Deployment

Models deployment in your IT environment.

## SCENARIO 2:

Data collection

Data Analysis

OCR **Dataset Creation** 

> "OCR" Optical Character Recognition is a powerful technological process that converts a paper document, PDF or even a photo into an editable digital document. Automated data extraction

The program behind this

technology starts by analyzing the structure of the image file, it tries to separate the different structures that compose the image: tables, texts, images, etc. Then each line is converted into words and then into

characters. Once the characters have been defined, the OCR begins to compare them with a group of images already processed by the OCR or predefined which allow it to propose hypotheses on the meaning of these characters. The OCR uses these

assumptions to analyze the curvatures of lines in words and words in characters. After reviewing all the assumptions, the program ultimately proposes content that could be consistent with the initial

Format error detection

Format errors, in the best-case scenario, can break an automated data processing

downstream analytical tasks that are difficult to debug. Example: international different way to write dates.

pipeline. In the worst case,

they introduce logical errors in

In our case, the documents are received in different formats.

Regex method assessment for SDS

**→** 

Sheet Regular Expression – better known as Regex – is pattern matching. They are used to match well-defined patterns (= matching tools identify words and phrases with common patterns). It is not an Al solutions using Natural

=> Focus on the Safety Data

Language Processing (NLP) and Machine Learning (ML) that can conduct an intelligent analysis of unstructured text. Indeed, NLP understands the language as it allows computers to undertake linguistic analysis and in practical terms "read" a document just like a human. In our situation, we don't need to understand the text. We need to report the content of document which are always

following the same pattern.

Apply the Regex model on SDS

on all SDS documents.

Apply the SDS regex model

=> Focus on the Test Report Regular Expression – better known as Regex – is pattern matching. They are used to match well-defined patterns (= matching tools identify words and phrases with common

Regex method

assessment for

TD

patterns). It is not an Al solutions using Natural Language Processing (NLP) and Machine Learning (ML) that can conduct an intelligent analysis of unstructured text. Indeed, NLP understands the language as it allows computers to undertake linguistic analysis and in practical terms "read" a document just like a human. In our situation, we don't need to understand the text. We need to report the content of document which are always

following the same pattern.

**ML Architecture** Design & Data labelling

The data labeling is the key step that will allows us to identify the screened data and provide one or more meaningful and informative labels. This will help the Machine Learning model in its learning process.

The organization of this algorithm will be personalised for the purpose of Bluesign needs and expectations.

ML training for TR Evaluation metrics reporting Model documentation

This is the final step =

machine learning model

settings. Combined, and

archive and current dataset),

Test Report document.

we will apply our model on the

Deployment

Models deployment in your IT combinaison of the raw data environment. analysis, features selection, applyed to the training set of the data (coming from your

SCENARIO 3:

**Dataset Creation** 

Data collection "OCR" Optical Character Data Analysis Recognition is a powerful technological process that converts a paper document, PDF or even a photo into an editable digital document. Automated data extraction

OCR (

Format errors, in the best-case scenario, can break an automated data processing pipeline. In the worst case, they introduce logical errors in downstream analytical tasks that are difficult to debug.

Format error

detection

(optional)

Regex method assessment for SDS

> Regular Expression – better known as Regex – is pattern matching. They are used to match well-defined patterns (=

Sheet

=> Focus on the Safety Data

Design & Data labelling

> step that will allows us to identify the screened data and provide one or more meaningful and informative labels. This will help the Machine Learning model in its

The data labeling is the key

**ML Architecture** 

SDS Evaluation metrics reporting Model documentation

ML training for

This is the final step = combinaison of the raw data analysis, features selection, machine learning model settings. Combined, and applyed to the training set of the data (coming from your

ML training for TR Evaluation metrics reporting Model documentation

Same thing than the previous step, but for Test Report instead of Safety Data Sheet.

Deployment

Models deployment in your IT environment.