# Data Extraction from Unstructured PDFs

This article was published as a part of the [Data Science Blogathon](#)

## Introduction:

Data Extraction is the process of extracting data from various sources such as CSV files, web, PDF, etc. Although in some files, data can be extracted easily as in CSV, while in files like unstructured PDFs we have to perform additional tasks to extract data.

There are a couple of Python libraries using which you can extract data from PDFs. For example, you can use the *PyPDF2* library for extracting text from PDFs where text is in a sequential or formatted manner i.e. in lines or forms. You can also extract tables in PDFs through the *Camelot* library. In all these cases data is in structured form i.e. sequential, forms or tables.

However, in the real world, most of the data is not present in any of the forms & there is no order of data. It is present in unstructured form. In this case, it is not feasible to use the above python libraries since they will give ambiguous results. To analyze unstructured data, we need to convert it to a structured form.

As such, there is no specific technique or procedure for extracting data from unstructured PDFs since data is stored randomly & it depends on what type of data you want to extract from PDF.

Here, I will show you a most successful technique & a python library through which you can extract data from *bounding boxes* in unstructured PDFs and then performing data cleaning operation on extracted data and converting it to a structured form.

## PyMuPDF:

I have used the PyMuPDF library for this purpose. This library provided many applications such as extracting images from PDF, extracting texts from different shapes, making annotations, draw a bounded box around the texts along with the features of libraries like *PyPDF2*.

Now, I will show you how I extracted data from the bounding boxes in a PDF with several pages.
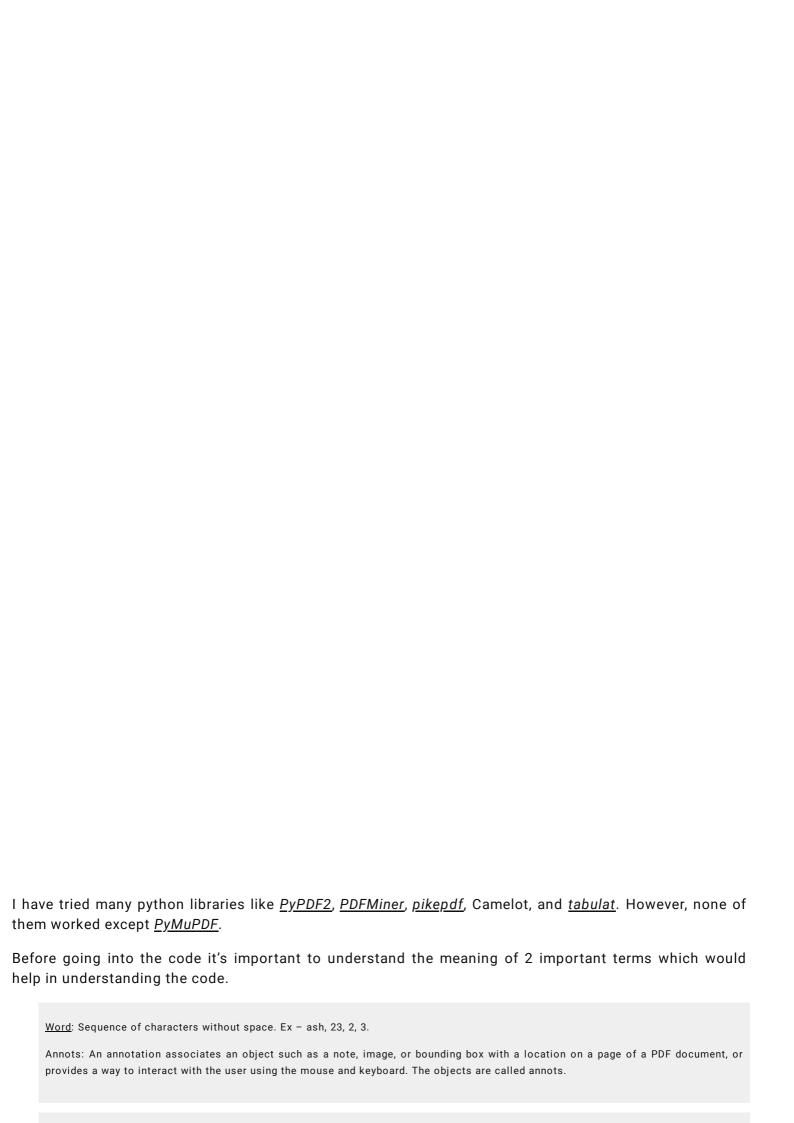
Here are the PDF and the red bounding boxes from which we need to extract data.

# TRAFFIC CRASH REPORT

*DENOTES MANDATORY FIELD FOR SUPPLEMENT REPORT

**LOCAL REPORT NUMBER ***  2 1 0 0 9 0 4 8

☑ PHOTOS TAKEN  ☑ OH-2  ☐ OH-3
☐ SECONDARY CRASH  ☐ OH-1P  ☐ OTHER
☐ PRIVATE PROPERTY

**LOCAL INFORMATION**
3F

**REPORTING AGENCY NAME *** MANSFIELD POLICE

**NCIC *** 0 7 0 0 1

**HIT/SKIP**
1 - SOLVED
2 - UNSOLVED

**NUMBER OF UNITS** 0 1

**UNIT IN ERROR**
98 - ANIMAL
99 - UNKNOWN
9 8

**COUNTY *** 7 0

**LOCALITY ***
1. CITY
2. VILLAGE
3. TOWNSHIP
1

**LOCATION: CITY, VILLAGE, TOWNSHIP *** MANSFIELD

**CRASH DATE/TIME *** 0 4 2 5 2 0 2 1 0 5 2 8

**CRASH SEVERITY**
1 - FATAL
2 - SERIOUS INJURY SUSPECTED
3 - MINOR INJURY SUSPECTED
4 - INJURY POSSIBLE
5 - PROPERTY DAMAGE ONLY
5

| ROUTE TYPE | ROUTE NUMBER | PREFIX 1-NORTH 2-SOUTH 3-EAST 4-WEST | LOCATION ROAD NAME | ROAD TYPE |
|---|---|---|---|---|
| | | 1 | MAIN | ST |

**LATITUDE** 4 0 . 7 9 0 6 2 0

| ROUTE TYPE | ROUTE NUMBER | PREFIX 1-NORTH 2-SOUTH 3-EAST 4-WEST | REFERENCE ROAD NAME(ROAD, MILEPOST, HOUSE #) | ROAD TYPE |
|---|---|---|---|---|
| | | | 1150 | |

**LONGITUDE** 8 2 . 5 1 3 2 1 0

**REFERENCE POINT**
1 - INTERSECTION
2 - MILE POST
3 - HOUSE #
3

**DIRECTION FROM REFERENCE**
1 - NORTH
2 - SOUTH
3 - EAST
4 - WEST

**ROUTE TYPE**
IR - INTERSTATE ROUTE (TP)
US - FEDERAL US ROUTE
SR - STATE ROUTE
CR - NUMBERED COUNTY ROUTE
TR - NUMBERED TOWNSHIP ROUTE

**ROAD TYPE**
AL - ALLEY
AV - AVENUE
BL - BOULEVARD
CR - CIRCLE
CT - COURT
DR - DRIVE
HE - HEIGHTS
HW - HIGHWAY
LA - LANE
MP - MILEPOST
OV - OVAL
PK - PARKWAY
PI - PIKE
PL - PLACE
RD - ROAD
SQ - SQUARE
ST - STREET
TE - TERRACE
TL - TRAIL
WA - WAY

**INTERSECTION RELATED**
☐ WITHIN INTERSECTION OR ON APPROACH
☐ WITHIN INTERCHANGE AREA
NUMBER OF APPROACHES

**DISTANCE FROM REFERENCE**

**DISTANCE UNIT OF MEASURE**
1 - MILES
2 - FEET
3 - YARDS

**ROADWAY**
☐ ROADWAY DIVIDED

**LOCATION OF FIRST HARMFUL EVENT**
1 - ON ROADWAY
2 - ON SHOULDER
3 - IN MEDIAN
4 - ON ROADSIDE
5 - ON GORE
6 - OUTSIDE TRAFFIC WAY
7 - ON RAMP
8 - OFF RAMP
9 - CROSSOVER
10 - DRIVEWAY/ALLEY ACCESS
11 - RAILWAY GRADE CROSSING
12 - SHARED USE PATHS OR TRAILS
13 - BIKE LANE
14 - TOLL BOOTH
99 - OTHER / UNKNOWN
1

**MANNER OF CRASH COLLISION/IMPACT**
1 - NOT COLLISION BETWEEN TWO MOTOR VEHICLES IN TRANSPORT
2 - REAR-END
3 - HEAD-ON
4 - REAR-TO-REAR
5 - BACKING
6 - ANGLE
7 - SIDESWIPE, SAME DIRECTION
8 - SIDESWIPE, OPPOSITE DIRECTION
9 - OTHER / UNKNOWN
1

**DIRECTION OF TRAVEL**
1 - NORTH
2 - SOUTH
3 - EAST
4 - WEST

**MEDIAN TYPE**
1 - DIVIDED FLUSH MEDIAN (<4 FEET)
2 - DIVIDED FLUSH MEDIAN (≥ 4 FEET)
3 - DIVIDED, DEPRESSED MEDIAN
4 - DIVIDED, RAISED MEDIAN (ANY TYPE)
9 - OTHER/UNKNOWN

☐ WORK ZONE RELATED
☐ WORKERS PRESENT
☐ LAW ENFORCEMENT PRESENT
☐ ACTIVE SCHOOL ZONE

**WORK ZONE TYPE**
1 - LANE CLOSURE
2 - LANE SHIFT/CROSSOVER
3 - WORK ON SHOULDER OR MEDIAN
4 - INTERMITTENT or MOVING WORK
5 - OTHER

**LOCATION OF CRASH IN WORK ZONE**
1 - BEFORE THE 1st WORK ZONE WARNING SIGN
2 - ADVANCE WARNING AREA
3 - TRANSITION AREA
4 - ACTIVITY AREA
5 - TERMINATION AREA

**CONTOUR** 1
1 - STRAIGHT LEVEL
2 - STRAIGHT GRADE
3 - CURVE LEVEL
4 - CURVE GRADE
9 - OTHER/UNKNOWN

**CONDITIONS** 1
1 - DRY
2 - WET
3 - SNOW
4 - ICE
5 - SAND, MUD, DIRT, OIL, GRAVEL
6 - WATER (STANDING, MOVING)
7 - SLUSH
9 - OTHER/UNKNOWN

**SURFACE** 2
1 - CONCRETE
2 - BLACKTOP, BITUMINOUS, ASPHALT
3 - BRICK/BLOCK
4 - SLAG, GRAVEL, STONE
5 - DIRT
9 - OTHER/UNKNOWN

**LIGHT CONDITIONS**
1 - DAYLIGHT
2 - DAWN/DUSK
3 - DARK - LIGHTED ROADWAY
4 - DARK - ROADWAY NOT LIGHTED
5 - DARK - UNKNOWN ROADWAY LIGHTING
9 - OTHER / UNKNOWN
3

**WEATHER**
1 - CLEAR
2 - CLOUDY
3 - FOG, SMOG, SMOKE
4 - RAIN
5 - SLEET, HAIL
6 - SNOW
7 - SEVERE CROSSWINDS
8 - BLOWING SAND, SOIL, DIRT, SNOW
9 - FREEZING RAIN OR FREEZING DRIZZLE
99 - OTHER/UNKNOWN
2

**NARRATIVE**

U#1 was Northbound on N Main St at 1150. A deer crossed the road Eastbound and U#1 struck it causing damage.

N MAIN ST    Block 1150

| CRASH REPORTED DATE/TIME | DISPATCH DATE/TIME | ARRIVAL TIME | SCENE CLEARED DATE / TIME |
|---|---|---|---|
| 0 4 2 5 2 0 2 1 0 5 2 8 | 0 4 2 5 2 0 2 1 0 5 3 2 | 0 4 2 5 2 0 2 1 0 5 3 6 | 0 4 2 5 2 0 2 1 0 5 4 7 |

**REPORT TAKEN BY**
☑ POLICE AGENCY
☐ MOTORIST
☐ SUPPLEMENT (CORRECTION OR ADDITION TO AN EXISTING REPORT SENT TO ODPS)

| TOTAL TIME ROADWAY CLOSED | OTHER INVESTIGATION TIME | TOTAL MINUTES | OFFICER'S NAME* | CHECKED BY OFFICER'S NAME* |
|---|---|---|---|---|
| | | 1 5 | CAROLYN YOUNG | P. WILLIAMS |

**OFFICER'S BADGE NUMBER*** 1 7 5 4

**CHECKED BY OFFICER'S BADGE NUMBER*** 1 0 5 1

HSY7001 OH1 1/19 [760-0820]

# MOTORIST / NON-MOTORIST

| UNIT # | NAME: LAST, FIRST, MIDDLE | DATE OF BIRTH | AGE | GENDER |
|---|---|---|---|---|
| 1 | STANTO, QUINN NICOLE | 0 4 1 3 1 9 8 0 | 4 1 | F |

**ADDRESS: CITY, STATE, ZIP**
1754 BROWNSTONE BLVD - H, TOLEDO, OH 43601

**CONTACT PHONE** - INCLUDE AREA CODE

| INJURIES | INJURED TAKEN BY | EMS AGENCY (NAME) | INJURED TAKEN TO: MEDICAL FACILITY (NAME, CITY) | SAFETY EQUIPMENT USED | DOT-C COMPLIANT MC HELMET | SEATING POSITION | AIR BAG USAGE | EJECTION | TRAPPED |
|---|---|---|---|---|---|---|---|---|---|
| 5 | | | | 4 | | 1 | 1 | 1 | 1 |

| OL STATE | OPERATOR LICENSE NUMBER | OFFENSE CHARGED | LOCAL CODE | OFFENSE DESCRIPTION | CITATION NUMBER |
|---|---|---|---|---|---|
| O H | TK825051 | | | | |

| OL CLASS | ENDORSEMENT SELECT UP TO 2 | RESTRICTION SELECT UP TO 3 | DRIVER DISTRACTED BY | ALCOHOL/DRUG SUSPECTED | CONDITION | ALCOHOL TEST | | | DRUG TEST(S) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | ☐ ALCOHOL ☐ MARIJUANA | | STATUS | TYPE | VALUE | STATUS | TYPE | RESULTS SELECT UP TO 4 |
| 4 | | | 1 | ☐ OTHER DRUG | 1 | 1 | | | 1 | | |

| UNIT # | NAME: LAST, FIRST, MIDDLE | DATE OF BIRTH | AGE | GENDER |
|---|---|---|---|---|
| | | | | |

**ADDRESS: CITY, STATE, ZIP**

**CONTACT PHONE** - INCLUDE AREA CODE

| INJURIES | INJURED TAKEN BY | EMS AGENCY (NAME) | INJURED TAKEN TO: MEDICAL FACILITY (NAME, CITY) | SAFETY EQUIPMENT USED | DOT-C COMPLIANT MC HELMET | SEATING POSITION | AIR BAG USAGE | EJECTION | TRAPPED |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | |

| OL STATE | OPERATOR LICENSE NUMBER | OFFENSE CHARGED | LOCAL CODE | OFFENSE DESCRIPTION | CITATION NUMBER |
|---|---|---|---|---|---|
| | | | | | |

| OL CLASS | ENDORSEMENT SELECT UP TO 2 | RESTRICTION SELECT UP TO 3 | DRIVER DISTRACTED BY | ALCOHOL/DRUG SUSPECTED | CONDITION | ALCOHOL TEST | | | DRUG TEST(S) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | ☐ ALCOHOL ☐ MARIJUANA | | STATUS | TYPE | VALUE | STATUS | TYPE | RESULTS SELECT UP TO 4 |
| | | | | ☐ OTHER DRUG | | | | | | | |

| UNIT # | NAME: LAST, FIRST, MIDDLE | DATE OF BIRTH | AGE | GENDER |
|---|---|---|---|---|
| | | | | |

**ADDRESS: CITY, STATE, ZIP**

**CONTACT PHONE** - INCLUDE AREA CODE

| INJURIES | INJURED TAKEN BY | EMS AGENCY (NAME) | INJURED TAKEN TO: MEDICAL FACILITY (NAME, CITY) | SAFETY EQUIPMENT USED | DOT-C COMPLIANT MC HELMET | SEATING POSITION | AIR BAG USAGE | EJECTION | TRAPPED |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | |

| OL STATE | OPERATOR LICENSE NUMBER | OFFENSE CHARGED | LOCAL CODE | OFFENSE DESCRIPTION | CITATION NUMBER |
|---|---|---|---|---|---|
| | | | | | |

| OL CLASS | ENDORSEMENT SELECT UP TO 2 | RESTRICTION SELECT UP TO 3 | DRIVER DISTRACTED BY | ALCOHOL/DRUG SUSPECTED | CONDITION | ALCOHOL TEST | | | DRUG TEST(S) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | ☐ ALCOHOL ☐ MARIJUANA | | STATUS | TYPE | VALUE | STATUS | TYPE | RESULTS SELECT UP TO 4 |
| | | | | ☐ OTHER DRUG | | | | | | | |

| INJURIES | SEATING POSITION | AIR BAG | OL CLASS | OL RESTRICTION(S) | DRIVER DISTRACTION | TEST STATUS |
|---|---|---|---|---|---|---|
| 1 - FATAL | 1 - FRONT – LEFT SIDE (MOTORCYCLE DRIVER) | 1 - NOT DEPLOYED | 1 - CLASS A | 1 - ALCOHOL INTERLOCK DEVICE | 1 - NOT DISTRACTED | 1 - NONE GIVEN |
| 2 - SUSPECTED SERIOUS INJURY | 2 - FRONT – MIDDLE | 2 - DEPLOYED FRONT | 2 - CLASS B | 2 - CDL INTRASTATE ONLY | 2 - MANUALLY OPERATING AN ELECTRONIC COMMUNICATION DEVICE (TEXTING, TYPING, DIALING) | 2 - TEST REFUSED |
| 3 - SUSPECTED MINOR INJURY | 3 - FRONT – RIGHT SIDE | 3 - DEPLOYED SIDE | 3 - CLASS C | 3 - CORRECTIVE LENSES | | 3 - TEST GIVEN, CONTAMINATED SAMPLE / UNUSABLE |
| 4 - POSSIBLE INJURY | 4 - SECOND – LEFT SIDE (MOTORCYCLE PASSENGER) | 4 - DEPLOYED BOTH FRONT / SIDE | 4 - REGULAR CLASS (OHIO = D) | 4 - FARM WAIVER | | 4 - TEST GIVEN, RESULTS KNOWN |
| 5 - NO APPARENT INJURY | 5 - SECOND – MIDDLE | 5 - NOT APPLICABLE | 5 - M/C MOPED ONLY | 5 - EXCEPT CLASS A BUS | 3 - TALKING ON HANDS-FREE COMMUNICATION DEVICE | 5 - TEST GIVEN, RESULTS UNKNOWN |
| | 6 - SECOND – RIGHT SIDE | 9 - DEPLOYMENT UNKNOWN | 6 - NO VALID OL | 6 - EXCEPT CLASS A & CLASS B BUS | | |
| **INJURED TAKEN BY** | 7 - THIRD – LEFT SIDE (MOTORCYCLE SIDE CAR) | | | 7 - EXCEPT TRACTOR-TRAILER | 4 - TALKING ON HAND-HELD COMMUNICATION DEVICE | **ALCOHOL TEST TYPE** |
| 1 - NOT TRANSPORTED / TREATED AT SCENE | 8 - THIRD – MIDDLE | **EJECTION** | **OL ENDORSEMENT** | 8 - INTERMEDIATE LICENSE RESTRICTIONS | 5 - OTHER ACTIVITY WITH AN ELECTRONIC DEVICE | 1 - NONE |
| 2 - EMS | 9 - THIRD – RIGHT SIDE | 1 - NOT EJECTED | H - HAZMAT | 9 - LEARNER'S PERMIT RESTRICTIONS | 6 - PASSENGER | 2 - BLOOD |
| 3 - POLICE | 10 - SLEEPER SECTION OF TRUCK CAB | 2 - PARTIALLY EJECTED | M - MOTORCYCLE | 10 - LIMITED TO DAYLIGHT ONLY | 7 - OTHER DISTRACTION INSIDE THE VEHICLE | 3 - URINE |
| 9 - OTHER / UNKNOWN | 11 - PASSENGER IN OTHER ENCLOSED CARGO AREA (NON-TRAILING UNIT, BUS, PICK-UP WITH CAP) | 3 - TOTALLY EJECTED | P - PASSENGER | 11 - LIMITED TO EMPLOYMENT | | 4 - BREATH |
| **SAFETY EQUIPMENT** | | 4 - NOT APPLICABLE | N - TANKER | 12 - LIMITED – OTHER | 8 - OTHER DISTRACTION OUTSIDE THE VEHICLE | 5 - OTHER |
| 1 - NONE USED | | | Q - MOTOR SCOOTER | 13 - MECHANICAL DEVICES (SPECIAL BRAKES, HAND CONTROLS, OR OTHER ADAPTIVE DEVICES) | 9 - OTHER / UNKNOWN | |
| 2 - SHOULDER BELT ONLY USED | 12 - PASSENGER IN UNENCLOSED CARGO AREA | **TRAPPED** | R - THREE-WHEEL MOTORCYCLE | | | **DRUG TEST TYPE** |
| 3 - LAP BELT ONLY USED | 13 - TRAILING UNIT | 1 - NOT TRAPPED | S - SCHOOL BUS | | **CONDITION** | 1 - NONE |
| 4 - SHOULDER & LAP BELT USED | 14 - RIDING ON VEHICLE EXTERIOR (NON-TRAILING UNIT) | 2 - EXTRICATED BY MECHANICAL MEANS | T - DOUBLE & TRIPLE TRAILERS | 14 - MILITARY VEHICLES ONLY | 1 - APPARENTLY NORMAL | 2 - BLOOD |
| 5 - CHILD RESTRAINT SYSTEM - FORWARD FACING | 15 - NON-MOTORIST | 3 - FREED BY NON-MECHANICAL MEANS | X - TANKER / HAZMAT | 15 - MOTOR VEHICLES WITHOUT AIR BRAKES | 2 - PHYSICAL IMPAIRMENT | 3 - URINE |
| 6 - CHILD RESTRAINT SYSTEM - REAR FACING | 99 - OTHER / UNKNOWN | | | 16 - OUTSIDE MIRROR | 3 - EMOTIONAL (E.G., DEPRESSED, ANGRY, DISTURBED) | 4 - OTHER |
| 7 - BOOSTER SEAT | | | **GENDER** | 17 - PROSTHETIC AID | 4 - ILLNESS | |
| 8 - HELMET USED | | | F - FEMALE | 18 - OTHER | 5 - FELL ASLEEP, FAINTED, FATIGUED, ETC. | **DRUG TEST RESULT(S)** |
| 9 - PROTECTIVE PADS USED (ELBOW, KNEES, ETC.) | | | M - MALE | | 6 - UNDER THE INFLUENCE OF MEDICATIONS / DRUGS / ALCOHOL | 1 - AMPHETAMINES |
| 10 - REFLECTIVE CLOTHING | | | U - OTHER/UNKNOWN | | 9 - OTHER / UNKNOWN | 2 - BARBITURATES |
| 11 - LIGHTING – PEDESTRIAN / BICYCLE ONLY | | | | | | 3 - BENZODIAZEPINES |
| 99 - OTHER / UNKNOWN | | | | | | 4 - CANNABINOIDS |
| | | | | | | 5 - COCAINE |
| | | | | | | 6 - OPIATES / OPIOIDS |
| | | | | | | 7 - OTHER |
| | | | | | | 8 - NEGATIVE RESULTS |

HSY8386 OH1M 1/19 [760-1500]

I have tried many python libraries like _PyPDF2_, _PDFMiner_, _pikepdf_, Camelot, and _tabulat_. However, none of them worked except _PyMuPDF_.

Before going into the code it's important to understand the meaning of 2 important terms which would help in understanding the code.

Word: Sequence of characters without space. Ex – ash, 23, 2, 3.

Annots: An annotation associates an object such as a note, image, or bounding box with a location on a page of a PDF document, or provides a way to interact with the user using the mouse and keyboard. The objects are called annots.

First, we will extract text from one of the bounding boxes. Then we will use the same procedure to extract data from all the bounding boxes of pdf.

# Code:

```
import fitz import pandas as pd doc = fitz.open('Mansfield--70-21009048 - ConvertToExcel.pdf') page1 = doc[0]
words = page1.get_text("words")
```

Firstly, we import the *fitz* module of the *PyMuPDF* library and pandas library. Then the object of the PDF file is created and stored in doc and 1st page of pdf is stored on page1. *page.get_text()* extracts all the words of page 1. Each word consists of a tuple with 8 elements.

In words variable, the First 4 elements represent the coordinates of the word, 5th element is the word itself, 6th,7th, 8th elements are block, line, word numbers respectively.

**OUTPUT**

**Extract the coordinates of the first object :**

```
first_annots=[] rec=page1.first_annot.rect rec #Information of words in first object is stored in mywords
mywords = [w for w in words if fitz.Rect(w[:4]) in rec] ann= make_text(mywords) first_annots.append(ann)
```

**This function selects the words contained in the box, sort the words and return in form of a string :**

```
def make_text(words): line_dict = {} words.sort(key=lambda w: w[0]) for w in words: y1 = round(w[3], 1) word
= w[4] line = line_dict.get(y1, []) line.append(word) line_dict[y1] = line lines = list(line_dict.items())
lines.sort() return "n".join([" ".join(line[1]) for line in lines])
```

**OUTPUT**

*page.first_annot()* gives the first annot i.e. bounding box of the page.

*.rect* gives coordinates of a rectangle.

Now, we got the coordinates of the rectangle and all the words on the page. We then filter the words which are present in our bounding box and store them in *mywords* variable.

We have got all the words in the rectangle with their coordinates. However, these words are in random order. Since we need the text sequentially and that only makes sense, we used a function make_text() which first sorts the words from left to right and then from top to bottom. It returns the text in string format.

Hurrah! We have extracted data from one annot. Our next task is to extract data from all annots of the PDF which would be done in the same approach.

### Extracting each page of the document and all the annots/rectanges :

```
for  pageno  in  range(0,len(doc)-1):  page  =  doc[pageno]  words  =  page.get_text("words")  for  annot  in
page.annots(): if annot!=None: rec=annot.rect mywords = [w for w in words if fitz.Rect(w[:4]) in rec] ann=
make_text(mywords) all_annots.append(ann)
```

*all_annots*, a list is initialized to store the text of all annots in the pdf.

The function of the outer loop in the above code is to go through each page of PDF, while that of the inner loop is to go through all annots of the page and performing the task of adding texts to all_annots list as discussed earlier.

Printing all_annots provides us the text of all annots of the pdf which you can see below.

### OUTPUT

Finally, we have extracted the texts from all the annots/ bounding boxes.

Its time to clean the data and bring it in an understandable form.

# Data Cleaning and Data Processing

### Splitting to form column name and its values :

```
cont=[] for i in range(0,len(all_annots)): cont.append(all_annots[i].split('n',1))
```

## Removing unnecessary symbols *,#,:

```
liss=[] for i in range(0,len(cont)): lis=[] for j in cont[i]: j=j.replace('*','') j=j.replace('#','')
j=j.replace(':','') j=j.strip() #print(j) lis.append(j) liss.append(lis)
```

## Spliting into keys and values and removing spaces in the values which only contain digits :

```
keys=[] values=[] for i in liss: keys.append(i[0]) values.append(i[1]) for i in range(0, len(values)): for j
in range(0,len(values[i])): if values[i][j]>='A' and values[i][j]<='Z': break if j==len(values[i])-1:
values[i]=values[i].replace(' ','')
```

We split each string based on a new line (n) character to separate the column name from its values. By further cleaning unnecessary symbols like (*, #, ⬚ are removed. Spaces between digits are removed.

With the key-value pairs, we create a dictionary which is shown below:

## Converting to dictionary :

```
report=dict(zip(keys,values))
```
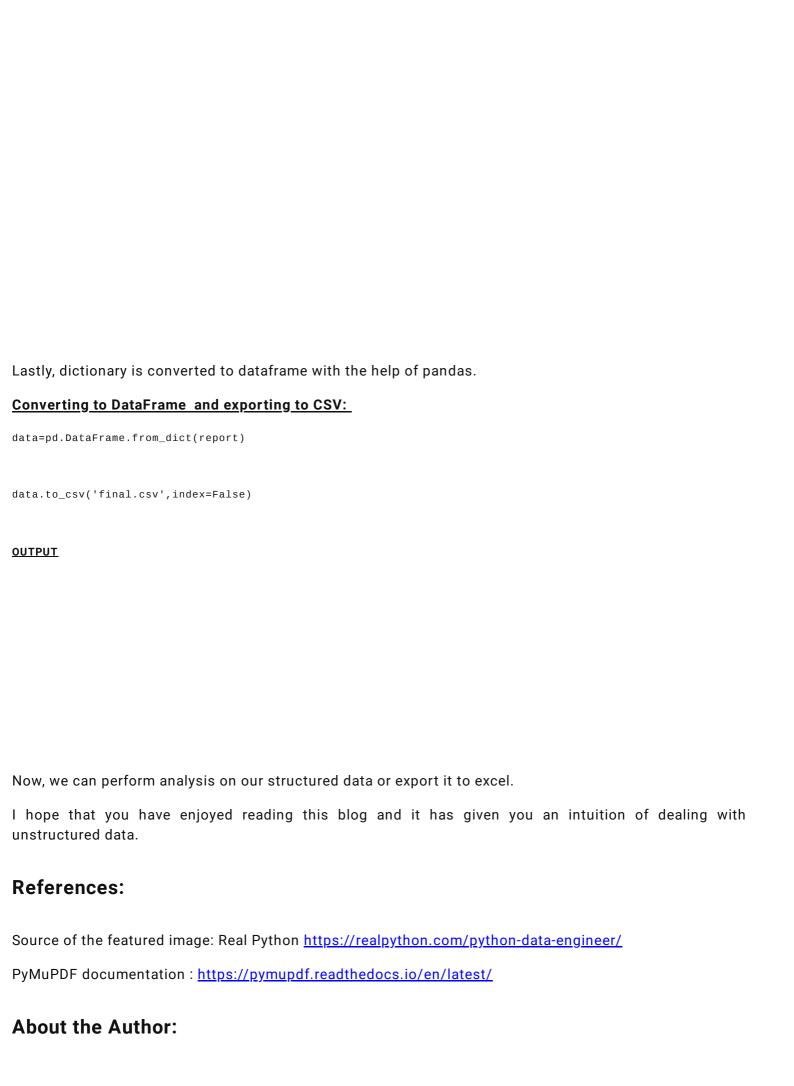
## report['VEHICLE IDENTIFICATION']=report['VEHICLE IDENTIFICATION'].replace(' ',")

```
dic=[report['LOCALITY'],report['MANNER OF CRASH COLLISION/IMPACT'],report['CRASH SEVERITY']] l=0 val_after=[]
for local in dic: li=[] lii=[] k='' extract='' l=0 for i in range(0,len(local)-1): if local[i+1]>='0' and
local[i+1]<='9': li.append(local[l:i+1]) l=i+1 li.append(local[l:]) print(li) for i in li: if i[0] in lii:
k=i[0] break lii.append(i[0]) for i in li: if i[0]==k:
```

## extract=i

```
val_after.append(extract)      break      report['LOCALITY']=val_after[0]      report['MANNER  OF  CRASH
COLLISION/IMPACT']=val_after[1] report['CRASH SEVERITY']=val_after[2]
```

## OUTPUT

Lastly, dictionary is converted to dataframe with the help of pandas.

**<u>Converting to DataFrame and exporting to CSV:</u>**

```
data=pd.DataFrame.from_dict(report)
```

```
data.to_csv('final.csv',index=False)
```

**<u>OUTPUT</u>**

Now, we can perform analysis on our structured data or export it to excel.

I hope that you have enjoyed reading this blog and it has given you an intuition of dealing with unstructured data.

# References:

Source of the featured image: Real Python [https://realpython.com/python-data-engineer/](https://realpython.com/python-data-engineer/)

PyMuPDF documentation : [https://pymupdf.readthedocs.io/en/latest/](https://pymupdf.readthedocs.io/en/latest/)

# About the Author:

Hi! I am Ashish Choudhary. I am pursuing B.Tech from the JC Bose University of Science & Technology. Data Science is my passion and feels proud to write interesting blogs related to it. Feel free to contact me on Linkedin linkedin.com/in/ashish-choudhary-7b6029166.

*The media shown in this article are not owned by Analytics Vidhya and are used at the Author's discretion.*

---

Article Url - https://www.analyticsvidhya.com/blog/2021/06/data-extraction-from-unstructured-pdfs/

**ashish@choudhary**