# FUZZY MATCHING: Conditional FUZZ PARTIAL TOKEN MATCH for Keyword2

## **Description:**

This is a Fuzzy Matching Project. There are two Files involved.

1. Raw data and

2. Keyword List

#### Raw data: Includes columns like:

1	1 raw_data.head(2)				
	HS_Code	Detailed_Product	HS_sub		
0	90181100	trolley for btl -08 ecg l line baik / baru	9018		
1	90183200	microcatheter.merit maestro2.8f.2.4f.130 cm	9018		

### **Keyword list looks like this:**

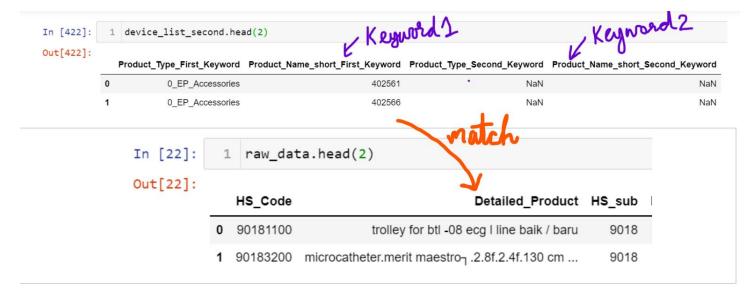
	Α	В	
1	Product_Name_short_First Keyword	Product_Name_short_Second Keyword	
2	dialysis	catheter	
3	Magnum	magnum ii	
4	Biopsy	plier	
5	Biopsy	nipper	
6	Biopsy	pincer	
7	Biopsy	forcep	
8	Coseal	CARBOSEAL	
9	Gelita	thp2628x150b	
10	Gelita	thp2830x100b	

The Keyword File has two main columns:

- Keyword 1 --> Product\_Name\_short\_First\_Keyword
- Keyword 2 --> Product\_Name\_short\_Second\_Keyword

# Step 1: Run the Fuzz Partial Match on raw\_data['Detailed\_Product'] for Keyword 1

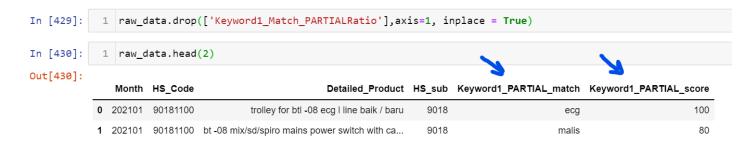
We import Keyword File and run Fuzz Partial Match for matching Keyword1, that is "Product\_Name\_short\_First\_Keyword" with the "raw\_data['Detailed\_Product']" and get the match score.



## **Running Fuzz Partial Token match for Keyword1**

```
In [390]:
            1 raw_data = raw_data.apply(lambda x: x.astype(str).str.lower())
In [425]:
            1 choices_brand = device_list_second['Product_Name_short_First_Keyword']
In [426]:
            1 choices_brand.shape
Out[426]: (9847,)
In [427]:
            1 def matchfunc_pr(x,choices):
                  option = process.extract(x, choices, limit=1, scorer=fuzz.partial_ratio)
                  if option[0][1]>75:
            4
                      return option
                      return ""
In [428]:
            1 raw_data['Keyword1_Match_PARTIALRatio']= raw_data['Detailed_Product'].apply(lambda x: matchfunc_pr(x,choices_brand))
            2 raw_data['Keyword1_PARTIAL_match'] = raw_data['Keyword1_Match_PARTIALRatio'].apply(lambda x: x[0][0] if x != '' else '')
            3 raw_data['Keyword1_PARTIAL_score'] = raw_data['Keyword1_Match_PARTIALRatio'].apply(lambda x: x[0][1] if x != '' else '')
```

#### And the Output generated looks like:



# Step 2 : Run the Fuzz Partial Match on raw\_data['Detailed\_Product'] for Keyword 2

Logic: If Keyword 1 is present in raw\_data [Keyword\_PARTIAL\_match] then select the corresponding Keyword 2 and run Fuzz Partial Token match for those Keywords, run Partial Token Match against "Detailed\_Product" and get the Score

1e_s	hort_Second Keyword 2, granysis	catneter
	16A	The B
1	Product_Name_short_First Keyword	Product_Name_short_Second Keyword
2	dialysis	catheter
3	Magnum	magnum ii
4	Biopsy	plier ]
5	Biopsy	nipper 🕊
6	Biopsy	pincer
7	Biopsy J	forcep
8	Coseal	CARBOSEAL
9	Gelita /	thp2628x150b
10	Gelita	thp2830x100b
11	Gelita	thp3032150b

### Approach:

1. Now our Keyword list would be raw\_data: Keyword1\_PARTIAL\_match

	Month	HS_Code	Detailed_Product	HS_sub	Keyword1_PARTIAL_match	Keyword1_PARTIAL_score
0	202101	90181100	trolley for btl -08 ecg I line baik / baru	9018	ecg	100
1	202101	90181100	bt -08 mix/sd/spiro mains power switch with ca	9018	malis	80

We copy that column to an empty List :

K1list = raw\_data['Keyword1\_PARTIAL\_match']

Keep only the unique values:

K1list = list(set(K1list))

#### - Then we create a while loop:

- For length in range equal to the len of Kllist
- If the first word in K1list matches "Product\_Name\_short\_First\_Keyword" then extract the adjacent values of

"Product\_Name\_short\_First\_Keyword" that is "Product\_Name\_short\_Second\_Keyword" to the

new column : \*\*device list second"['K2']\*\*

- Now these chunk of keywords in device\_list\_second"['K2'] will be the keyword s to be matched.
- So, we check for the raw\_data.loc where "Keyword1\_PARTIAL\_Match" is equal to first word in K1list
- We now have a slice of raw\_data (df\_temp)
- What the While loop is doing is , getting the Keywords2 in device\_list\_secon  $\tt d"['K2']$

- Matching those words against the "Detailed\_Product" where first word in Klli st matches the Keywordl PARTIAL Match
- We run Fuzz partial match for df temp
- Keep appending this newly updated slice  $(df\_temp)$  to the Dataframe list as a nd when the new
- df temp is generated in the While loop.
- Towards the end we have the Final df "data" with all the df temp appended.

```
#Fuzz partial match Function that we are applying to the sclice of data in the while loop
  2
     def matchfunc_pr(x,choices):
  3
           option = process.extract(x, choices, limit=1, scorer=fuzz.partial_ratio)
  4
  5
  6
           if option[0][1]>75:
  7
               return option
  8
           else:
                return ""
  9
In [ ]: 1 | K1list = raw_data['Keyword1_PARTIAL_match']
         2 K1list = list(set(K1list))
         3 K1list.remove("")
         4 data = pd.DataFrame([])
         6 i = 0
         7 while i < len(K1list):</pre>
         8
        9
                   device list second['K2'] = device list second.loc
                           [device_list_second['Product_Name_short_First_Keyword'] == K1list[i], 'Product_Name_short_Second_Keyword']
        10
                   device_list_second['K2'] = device_list_second['K2'].astype(str)
        11
        12
                   choices_brand =device_list_second['K2']
        13
                   df_temp = raw_data.loc[raw_data['Keyword1_PARTIAL_match'] == K1list[i]]
        14
                   df_temp = df_temp.copy()
                   df_temp['Keyword2_Match_PARTIALRatio'] = df_temp['Detailed_Product'].apply(lambda x: matchfunc_pr(x,choices_brand))
        15
                   df_temp['Keyword2_PARTIAL_match'] = df_temp['Keyword2_Match_PARTIALRatio'].
        16
                                                             apply(lambda x: x[0][0] if x != '' else '')
        17
                   df_temp['Keyword2_PARTIAL_score'] = df_temp['Keyword2_Match_PARTIALRatio'].
        18
        19
                                                             apply(lambda x: x[0][1] if x != '' else '')
                   data = data.append(df_temp, ignore_index=True)
        20
        21
                   i = i + 1
```

### We have the new Data Frame which looks like:

```
In [433]: 1 data.shape
Out[433]: (74441, 9)
In [434]: 1 data.head(2)
Out[434]:
            Month HS_Code Detailed_Product HS_sub Keyword1_PARTIAL_match Keyword1_PARTIAL_score Keyword2_Match_PARTIALRatio Keyword2_PARTIAL_match
                              blade f/bv380r 4-
            202109 90189090
                               prong 16mm x
                                                                                                83
                                      60mm
                              blade f/bv380r 4-
            202109 90189090
                               prong 16mm x
                                               9018
                                                                     bv580r
                                                                                                83
                                      40mm
```

The Column "Keyword2\_PARTIAL\_match" will contain the matched word.

The column "Keyword2\_PARTIAL\_score" will contain the match Score.

Owner : Aisha Khalid