# ddh & daeyou EDA

끝이나지 않는 탐색전

# 🔍 목차

# 과제내용



gt : 실질 치아 포인터 값

guess : 인공지능이 인식한 치아 포인터 값

x,y,z : 치아 포인터별 좌표

# 탐색이라고 하고 헤엄

| 데이터 전처리 | visualization | χ² - test | t-test |
|---|---|---|---|
| 데이터 merge | Boxen Plot | 범주형 data 상관관계분석 | 범주형 data -연속형 data 차이분석 |
| 데이터 타입 분류 | 2D-Scatter Plot | | |
| | KDE Plot | | |

# PROJECT 1

# data merge

```
#파일 불러오기
ddh=pd.read_csv(r'C:\Users\imagoworks-moongzeee\Desktop\과제\2021.01.21 과제\ddh-eda.csv')
daeyou=pd.read_csv(r'C:\Users\imagoworks-moongzeee\Desktop\과제\2021.01.21 과제\daeyou-eda.csv')

#ddh, daeyou dataframe
ddh.columns=['x','y','z','gt','guess']
daeyou.columns=['x','y','z','gt','guess']
ddh['category']='ddh'
daeyou['category']='daeyou'

#ddh와 daeyou의 gt&guess 열을 비교하기 위해 numpy로 변환
ddh_gt=np.array(ddh['gt'])
ddh_guess=np.array(ddh['guess'])
daeyou_gt=np.array(daeyou['gt'])
daeyou_guess=np.array(daeyou['guess'])

#ddh, daeyou의 gt와 guess를 비교해 boolean형태로 반환
ddh_A=pd.Series(ddh_gt==ddh_guess)
daeyou_A=pd.Series(daeyou_gt==daeyou_guess)

#gt와 guess 비교한 accuracy 열 DATA dataframe에 삽입
ddh['Accuracy']=ddh_A
daeyou['Accuracy']=daeyou_A

#merge
DATA=pd.concat([ddh,daeyou],ignore_index=True)
```

|   | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 |   | x | y | z | gt | guess | category | Accuracy |
| 2 | 0 | -24.35 | -6.68381 | -25.3205 | 0 | 0 | ddh | TRUE |
| 3 | 1 | -24.2311 | -6.90204 | -25.3443 | 0 | 0 | ddh | TRUE |
| 4 | 2 | -24.2656 | -6.80062 | -25.411 | 0 | 0 | ddh | TRUE |

# PROJECT 1

# data type change

```
8    DATA=pd.read_csv(r'C:\Users\imagoworks-moongzeee\Desktop\과제\2021.01.21 과제\EDA_DATA.csv')
9    print(DATA.info())
10   DATA['gt']=DATA['gt'].astype(object)
11   DATA['guess']=DATA['guess'].astype(object)
12   DATA['Accuracy']=DATA['Accuracy'].astype(object)
13   print(DATA.info())
14
15
16
```

PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL

```
(py) PS C:\Users\imagoworks-moongzeee\test> & C:/Users/imagoworks-moongzeee/.conda/envs/py/python.exe c:/Us
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 145082 entries, 0 to 145081
Data columns (total 8 columns):
 #   Column      Non-Null Count   Dtype
---  ------      --------------   -----
 0   Unnamed: 0  145082 non-null  int64
 1   x           145082 non-null  float64
 2   y           145082 non-null  float64
 3   z           145082 non-null  float64
 4   gt          145082 non-null  int64
 5   guess       145082 non-null  int64
 6   category    145082 non-null  object
 7   Accuracy    145082 non-null  bool
dtypes: bool(1), float64(3), int64(3), object(1)
memory usage: 7.9+ MB
None
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 145082 entries, 0 to 145081
Data columns (total 8 columns):
 #   Column      Non-Null Count   Dtype
---  ------      --------------   -----
 0   Unnamed: 0  145082 non-null  int64
 1   x           145082 non-null  float64
 2   y           145082 non-null  float64
 3   z           145082 non-null  float64
 4   gt          145082 non-null  object
 5   guess       145082 non-null  object
 6   category    145082 non-null  object
 7   Accuracy    145082 non-null  object
dtypes: float64(3), int64(1), object(4)
memory usage: 8.9+ MB
None
```

# Boxen Plot

## ddh

## daeyou

## 2D-Scatter Plot

### ddh





### daeyou

# KDE Plot

# 기술통계분석

## pandas profilereport module



보고서가 html 파일로 생성됨

# 문득 pointer가 뭔지 궁금해서



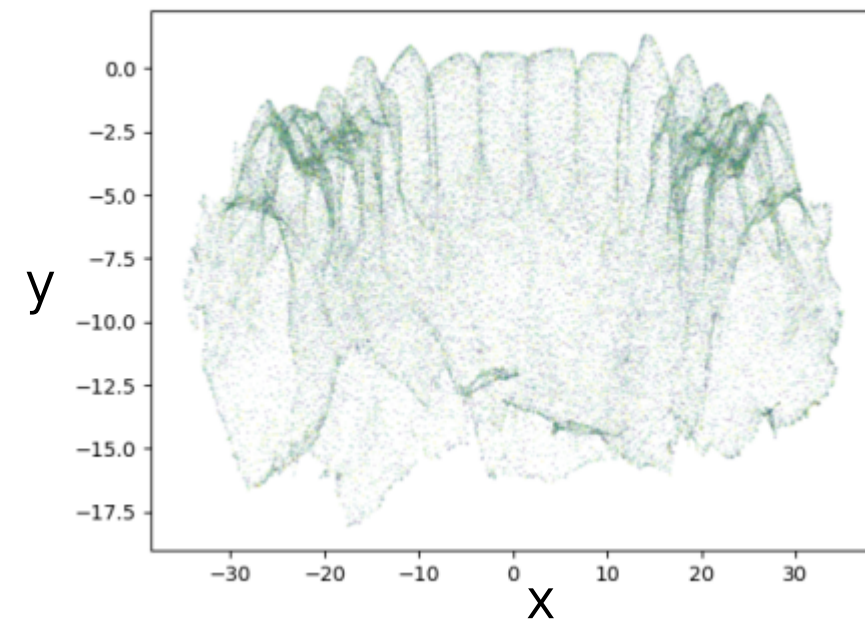| | ddh pointer별 갯수 | ddh pointer별 상대비율 | daeyou pointer별 갯수 | daeyou pointer별 상대비율 |
|---|---|---|---|---|
| 0 | 24343 | 0.472606 | 58942 | 0.629897 |
| 1 | 1260 | 0.024462 | 1555 | 0.016618 |
| 2 | 1404 | 0.027258 | 1770 | 0.018916 |
| 3 | 1642 | 0.031879 | 2050 | 0.021908 |
| 4 | 1712 | 0.033238 | 2458 | 0.026268 |
| 5 | 1566 | 0.030403 | 2378 | 0.025413 |
| 6 | 3207 | 0.062262 | 3977 | 0.042501 |
| 7 | 2964 | 0.057544 | 3347 | 0.035768 |
| 9 | 1193 | 0.023161 | 1525 | 0.016297 |
| 10 | 1234 | 0.023957 | 1788 | 0.019108 |
| 11 | 1496 | 0.029044 | 2252 | 0.024067 |
| 12 | 1594 | 0.030947 | 2206 | 0.023575 |
| 13 | 1589 | 0.030850 | 2210 | 0.023618 |
| 14 | 3361 | 0.065252 | 4021 | 0.042971 |
| 15 | 2943 | 0.057137 | 3095 | 0.033075 |



0 : 잇몸

# DDH vs DAEYOU

```
82    Accuracy=pd.DataFrame(
83        {'ddh count' :ddh_A.value_counts(normalize=False),
84         'ddh percent':ddh_A.value_counts(normalize=True),
85         'daeyou count':daeyou_A.value_counts(normalize=False),
86         'daeyou percent':daeyou_A.value_counts(normalize=True)}
87    )
88
89
90    print(Data_count)
91    #print(pointer_count)
92    print(Accuracy)
93    |
94
PROBLEMS  6    OUTPUT    DEBUG CONSOLE    TERMINAL

(py) PS C:\Users\imagoworks-moongzeee\test> & C:/Users/imagoworks-moongzeee
   ddh_data  daeyou_data
0     51508      93574
         ddh count   ddh percent  daeyou count  daeyou percent
True        49592      0.962802         90597        0.968186
False        1916      0.037198          2977        0.031814
(py) PS C:\Users\imagoworks-moongzeee\test> []
```

DDH Accuracy True의 비율

96.28%

DAEYOU Accuracy True 비율

96.81%

DAEYOU 가

**0.2 %** 더 정확함

# χ² - test

**H0 : 두 범주간의 관계가 독립적**

**H1 : 두 범주간의 관계가 종속적**

```python
#카이제곱 검정을 이용해 ddh/daeyou의guess 와 accuracy 범주 교차분석

result1=pd.crosstab(DATA_ddh.guess,DATA_ddh.Accuracy)
a=ss.chi2_contingency(observed=result1)

print('ddh_guess와 accurcy 간의 chi-Square:'+str(a[0]))
print('p-value:'+str(a[1]))
print('degree of freedom:'+str(a[2]))
print('expectation array:\n'+str(a[3]))

result2=pd.crosstab(DATA_daeyou.guess,DATA_daeyou.Accuracy)
a=ss.chi2_contingency(observed=result2)
print('daeyou_guess와 accurcy 간의 chi-Square:'+str(a[0]))
print('p-value:'+str(a[1]))
print('degree of freedom:'+str(a[2]))
print('expectation array:\n'+str(a[3]))
```

```
ddh_guess와 accurcy 간의 chi-Square:241.46531189059303
p-value:1.6673742817383728e-43
degree of freedom:14
expectation array:
[[  935.94152365 24225.05847635]
 [   43.85656597  1135.14343403]
 [   46.16284849  1194.83715151]
 [   56.87590277  1472.12409723]
 [   62.75320339  1624.24679661]
 [   56.05754446  1450.94245554]
 [  121.26582278  3138.73417722]
 [  106.57257125  2758.42742875]
 [   42.92661334  1111.07338666]
 [   44.7865186   1159.2134814 ]
 [   52.63531879  1362.36468121]
 [   58.40102508  1511.59897492]
 [   57.24788382  1481.75211618]
 [  123.64650151  3200.35349849]
 [  106.87015609  2766.12984391]]
```

?

```
daeyou_guess와 accurcy 간의 chi-Square:1853.4583168356069
p-value:0.0
degree of freedom:14
expectation array:
[[1.86301903e+03 5.66959810e+04]
 [4.72443734e+01 1.43775563e+03]
 [5.13802445e+01 1.56361976e+03]
 [6.74783273e+01 2.05352167e+03]
 [8.27174215e+01 2.51728258e+03]
 [7.41593498e+01 2.25684065e+03]
 [1.26875799e+02 3.86112420e+03]
 [1.02951375e+02 3.13304862e+03]
 [4.84533204e+01 1.47454668e+03]
 [5.52934148e+01 1.68270659e+03]
 [7.63227285e+01 2.32267727e+03]
 [7.28231453e+01 2.21617685e+03]
 [7.19959711e+01 2.19100403e+03]
 [1.26494026e+02 3.84950597e+03]
 [1.09791470e+02 3.34120853e+03]]
```

daeyou_guess와 accuracy간의 p-value 값이 0.05 보다 낮으므로 H0 가설은 기각

즉, daeyou_guess 와 accuracy는 연관이 있다.

# t - test



daeyou와 ddh의 y-pointer 값의 차이가 있다

# t - test

**H0 : 두 집단간의 평균의 차이가 없다**

**H1 : 두 집단간의 평균의 차이가 있다**

```python
47    ## daeyou/ddh pointer t-test
48
49    tg_ddh=DATA_ddh['x']
50    tg_daeyou=DATA_daeyou['x']
51
52
53    #등분산성 계산
54    print(ss.levene(tg_ddh,tg_daeyou))
55    print(ss.fligner(tg_ddh,tg_daeyou))
56    print(ss.bartlett(tg_ddh,tg_daeyou))
57
58
59    #t-test
60    print(ss.ttest_ind(tg_ddh, tg_daeyou, equal_var=True))
61
62
```

```
PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL

(py) PS C:\Users\imagoworks-moongzeee\test> & C:/Users/imagoworks-moongzeee/.c
LeveneResult(statistic=1157.4765639175257, pvalue=1.0617496113576694e-252)
FlignerResult(statistic=1820.436764555484, pvalue=0.0)
BartlettResult(statistic=752.3930289898982, pvalue=1.2107988824625322e-165)
Ttest_indResult(statistic=4.734551594659522, pvalue=2.197467710379726e-06)
```

```python
47    ## daeyou/ddh pointer t-test
48
49    tg_ddh=DATA_ddh['z']
50    tg_daeyou=DATA_daeyou['z']
51
52
53    #등분산성 계산
54    print(ss.levene(tg_ddh,tg_daeyou))
55    print(ss.fligner(tg_ddh,tg_daeyou))
56    print(ss.bartlett(tg_ddh,tg_daeyou))
57
58
59    #t-test
60    print(ss.ttest_ind(tg_ddh, tg_daeyou, equal_var=False))
61
62
```

```
PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL

(py) PS C:\Users\imagoworks-moongzeee\test> & C:/Users/imagoworks-moongzeee/.c
LeveneResult(statistic=2543.5567777598, pvalue=0.0)
FlignerResult(statistic=3183.720041844067, pvalue=0.0)
BartlettResult(statistic=1957.9232151890924, pvalue=0.0)
Ttest_indResult(statistic=35.369500524004636, pvalue=1.2117641456041473e-272)
```

```python
47    ## daeyou/ddh pointer t-test
48
49    tg_ddh=DATA_ddh['y']
50    tg_daeyou=DATA_daeyou['y']
51
52
53    #등분산성 계산
54    print(ss.levene(tg_ddh,tg_daeyou))
55    print(ss.fligner(tg_ddh,tg_daeyou))
56    print(ss.bartlett(tg_ddh,tg_daeyou))
57
58
59    #t-test
60    print(ss.ttest_ind(tg_ddh, tg_daeyou, equal_var=False))
61
62
```

```
PROBLEMS    OUTPUT    DEBUG CONSOLE    TERMINAL

(py) PS C:\Users\imagoworks-moongzeee\test> & C:/Users/imagoworks-moon
LeveneResult(statistic=10326.44535344914, pvalue=0.0)
FlignerResult(statistic=10052.892996492645, pvalue=0.0)
BartlettResult(statistic=12245.086414284988, pvalue=0.0)
Ttest_indResult(statistic=145.3727767407498, pvalue=0.0)
```

daeyou_y와 ddh_y간의 p-value 값이 0.05 보다 낮으므로 H0 가설은 기각

즉, daeyou_y 와 ddh_y 의 평균값 차이가 있다.

# plotly 패키지

https://plotly.com/python/3d-scatter-plots/