2020 데이터사이언스 프로젝트 최종 발표

# ▾ 국내 병원의 개/폐업 예측

팀원 : 김영민, 강민정

본 프로젝트는 병원에 대한 회계 데이터를 통해 국내 병원 개/폐업 예측한다.

(1) KDD 분석 절차에 따라 진행되며 EDA의 시각화를 통해 데이터의 특징 파악

(2) 머신러닝 모델과 인공 신경망 모델의 성능을 비교하고 이를 통해 보다 높은 정확도를 얻어냄

```
1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import warnings
5 warnings.filterwarnings('ignore')
6 import seaborn as sns
7 plt.style.use('ggplot')
```

```
1 # 구글 드라이브 마운트
2 from google.colab import drive
3 drive.mount('/content/gdrive')
```

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=947318989803

Enter your authorization code:
· · · · · · · · · · ·
Mounted at /content/gdrive

# ▾ 데이터 설명

train.csv - 의료기관이 폐업했는지 여부를 포함하여 최근 2개년의 재무정보와 병원 기본정보

test.csv - 폐업 여부를 제외하고 train.csv와 동일

sample_submission.csv - inst_id와 open과 close를 예측하는 OC 두개의 열로 구성. OC의 값은 open 예측일 경우 1, close 예측일 경우 0.

inst_id - 각 파일에서의 병원 고유 번호

OC – 영업/폐업 분류, 2018년 폐업은 2017년 폐업으로 간주함

sido – 병원의 광역 지역 정보

sgg – 병원의 시군구 자료

openDate – 병원 설립일

bedCount - 병원이 갖추고 있는 병상의 수

instkind – 병원, 의원, 요양병원, 한의원, 종합병원 등 병원의 종류

· 종합병원 : 입원환자 100명 이상 수용 가능

· 병원 : 입원 환자 30명 이상 100명 미만 수용 가능

· 의원 : 입원 환자 30명 이하 수용 가능

· 한방 병원(한의원) : 침술과 한약으로 치료하는 의료 기관.

revenue1 – 매출액, 2017(회계년도)년 데이터를 의미함

salescost1 – 매출원가, 2017(회계년도)년 데이터를 의미함

sga1 - 판매비와 관리비, 2017(회계년도)년 데이터를 의미함

salary1 – 급여, 2017(회계년도)년 데이터를 의미함

noi1 – 영업외수익, 2017(회계년도)년 데이터를 의미함

noe1 – 영업외비용, 2017(회계년도)년 데이터를 의미함

Interest1 – 이자비용, 2017(회계년도)년 데이터를 의미함

ctax1 – 법인세비용, 2017(회계년도)년 데이터를 의미함

Profit1 – 당기순이익, 2017(회계년도)년 데이터를 의미함

liquidAsset1 – 유동자산, 2017(회계년도)년 데이터를 의미함

quickAsset1 – 당좌자산, 2017(회계년도)년 데이터를 의미함

receivableS1 - 미수금(단기), 2017(회계년도)년 데이터를 의미함

inventoryAsset1 – 재고자산, 2017(회계년도)년 데이터를 의미함

nonCAsset1 – 비유동자산, 2017(회계년도)년 데이터를 의미함

tanAsset1 – 유형자산, 2017(회계년도)년 데이터를 의미함

OnonCAsset1 - 기타 비유동자산, 2017(회계년도)년 데이터를 의미함

receivableL1 – 장기미수금, 2017(회계년도)년 데이터를 의미함

debt1 – 부채총계, 2017(회계년도)년 데이터를 의미함

liquidLiabilities1 – 유동부채, 2017(회계년도)년 데이터를 의미함

shortLoan1 – 단기차입금, 2017(회계년도)년 데이터를 의미함

NCLiabilities1 – 비유동부채, 2017(회계년도)년 데이터를 의미함

longLoan1 – 장기차입금, 2017(회계년도)년 데이터를 의미함

netAsset1 – 순자산총계, 2017(회계년도)년 데이터를 의미함

surplus1 – 이익잉여금, 2017(회계년도)년 데이터를 의미함

revenue2 – 매출액, 2016(회계년도)년 데이터를 의미함

salescost2 – 매출원가, 2016(회계년도)년 데이터를 의미함

sga2 - 판매비와 관리비, 2016(회계년도)년 데이터를 의미함

salary2 – 급여, 2016(회계년도)년 데이터를 의미함

noi2 – 영업외수익, 2016(회계년도)년 데이터를 의미함

noe2 – 영업외비용, 2016(회계년도)년 데이터를 의미함

interest2 – 이자비용, 2016(회계년도)년 데이터를 의미함

ctax2 – 법인세비용, 2016(회계년도)년 데이터를 의미함

profit2 – 당기순이익, 2016(회계년도)년 데이터를 의미함

liquidAsset2 – 유동자산, 2016(회계년도)년 데이터를 의미함

quickAsset2 – 당좌자산, 2016(회계년도)년 데이터를 의미함

receivableS2 - 미수금(단기), 2016(회계년도)년 데이터를 의미함

inventoryAsset2 – 재고자산, 2016(회계년도)년 데이터를 의미함

nonCAsset2 – 비유동자산, 2016(회계년도)년 데이터를 의미함

tanAsset2 – 유형자산, 2016(회계년도)년 데이터를 의미함

OnonCAsset2 - 기타 비유동자산, 2016(회계년도)년 데이터를 의미함

receivableL2 – 장기미수금, 2016(회계년도)년 데이터를 의미함

Debt2 – 부채총계, 2016(회계년도)년 데이터를 의미함

liquidLiabilities2 – 유동부채, 2016(회계년도)년 데이터를 의미함

shortLoan2 – 단기차입금, 2016(회계년도)년 데이터를 의미함

NCLiabilities2 – 비유동부채, 2016(회계년도)년 데이터를 의미함

longLoan2 – 장기차입금, 2016(회계년도)년 데이터를 의미함

netAsset2 – 순자산총계, 2016(회계년도)년 데이터를 의미함

surplus2 – 이익잉여금, 2016(회계년도)년 데이터를 의미함

employee1 – 고용한 총 직원의 수, 2017(회계년도)년 데이터를 의미함

employee2 – 고용한 총 직원의 수, 2016(회계년도)년 데이터를 의미함

ownerChange – 대표자의 변동

```
1 train = pd.read_csv('/content/gdrive/My Drive/hospital close or open/train.csv',parse_dates=['op
2 test=pd.read_csv('/content/gdrive/My Drive/hospital close or open/test.csv',parse_dates=['openDa
3 answer = pd.read_csv('/content/gdrive/My Drive/hospital close or open/submission_sample.csv')
```

```
1 train_len = train.shape[0]
2 train_id = train.inst_id
3 test_len = test.shape[0]
4 test_id = test.inst_id
5 # 나중에 train과 test 따로 분류할 때 필요한 것들
```

```
6 merge_data = pd.concat([train,test]).reset_index() # data 병합
7 data=merge_data.copy()
```

```
1 test.shape
```

(127, 58)

## EDA

```
1 data.head()
```

| | index | inst_id | OC | sido | sgg | openDate | bedCount | instkind | re |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 1 | open | choongnam | 73 | 2007-12-28 | 175.0 | nursing_hospital | 4.2175 |
| **1** | 1 | 3 | open | gyeongnam | 32 | 1997-04-01 | 410.0 | general_hospital | |
| **2** | 2 | 4 | open | gyeonggi | 89 | 2016-12-28 | 468.0 | nursing_hospital | 1.0045 |
| **3** | 3 | 7 | open | incheon | 141 | 2000-08-14 | 353.0 | general_hospital | 7.2507 |
| **4** | 4 | 9 | open | gyeongnam | 32 | 2005-09-01 | 196.0 | general_hospital | 4.9043 |

```
1 train.shape
```

(301, 58)

```
1 # 결측치 확인
2 train.isnull().sum()/train.shape[0]*100
```
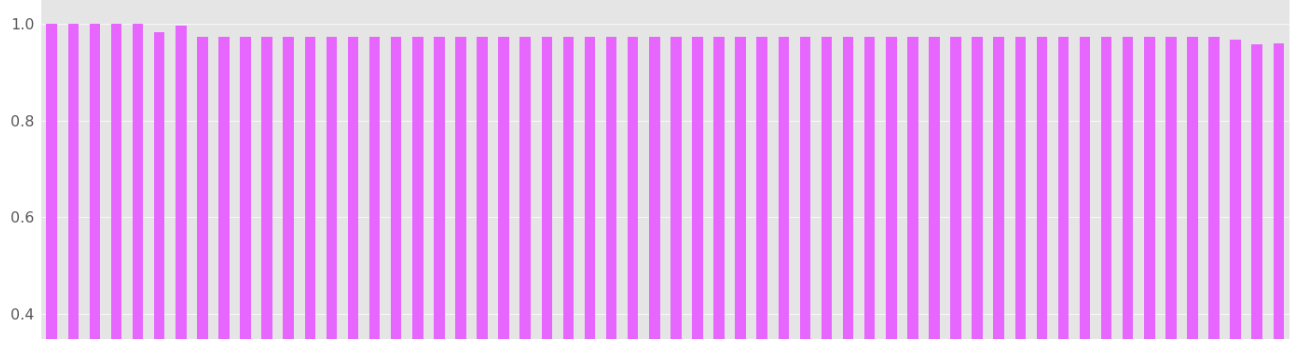
```
inst_id               0.000000
OC                    0.000000
sido                  0.000000
sgg                   0.000000
openDate              0.000000
bedCount              1.661130
instkind              0.332226
revenue1              2.657807
salescost1            2.657807
sga1                  2.657807
salary1               2.657807
noi1                  2.657807
noe1                  2.657807
interest1             2.657807
ctax1                 2.657807
profit1               2.657807
liquidAsset1          2.657807
quickAsset1           2.657807
receivableS1          2.657807
inventoryAsset1       2.657807
nonCAsset1            2.657807
tanAsset1             2.657807
OnonCAsset1           2.657807
receivableL1          2.657807
debt1                 2.657807
liquidLiabilities1    2.657807
shortLoan1            2.657807
NCLiabilities1        2.657807
longLoan1             2.657807
netAsset1             2.657807
surplus1              2.657807
revenue2              2.657807
salescost2            2.657807
sga2                  2.657807
salary2               2.657807
noi2                  2.657807
noe2                  2.657807
interest2             2.657807
ctax2                 2.657807
profit2               2.657807
liquidAsset2          2.657807
quickAsset2           2.657807
```

```
1 import missingno as msno
2 msno.bar(train,color=(0.9,0.4,1)) # 결측치가 거의 없음을 알 수 있다.
```

<matplotlib.axes._subplots.AxesSubplot at 0x7fdcfc24fba8>



```
1 data.describe(include='object')
```

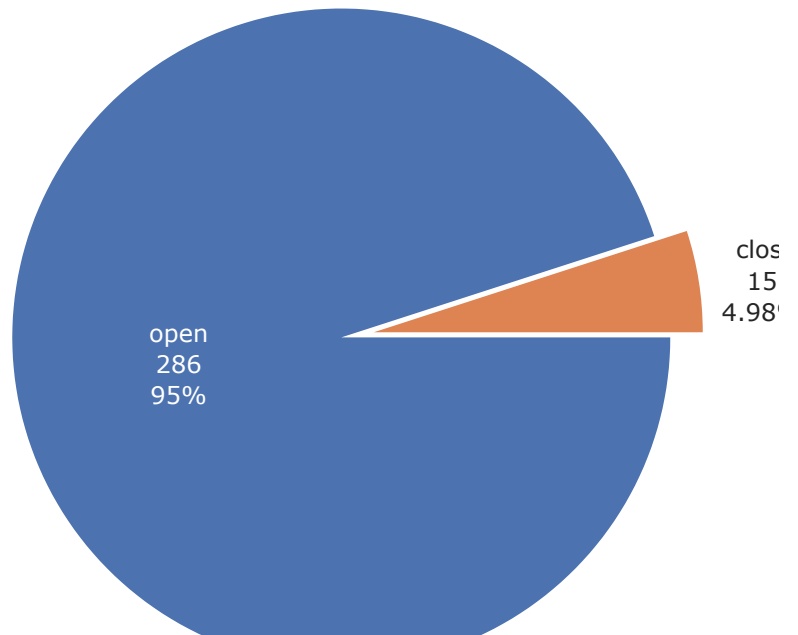|       | OC    | sido     | instkind         | employee1 | employee2 | ownerChange |
|-------|-------|----------|------------------|-----------|-----------|-------------|
| count | 301   | 428      | 425              | 410.0     | 400.0     | 401         |
| unique| 2     | 17       | 7                | 268.0     | 263.0     | 2           |
| top   | open  | gyeonggi | nursing_hospital | 73.0      | 70.0      | same        |
| freq  | 286   | 68       | 207              | 6.0       | 6.0       | 348         |

```
1 data.describe()
```

|       | index      | inst_id    | sgg        | bedCount   | revenue1     | salescost1   |        |
|-------|------------|------------|------------|------------|--------------|--------------|--------|
| count | 428.000000 | 428.000000 | 428.000000 | 415.000000 | 4.180000e+02 | 4.180000e+02 | 4.180  |
| mean  | 124.184579 | 215.154206 | 85.436916  | 153.474699 | 1.448543e+10 | 2.600554e+09 | 1.125  |
| std   | 85.469723  | 124.453370 | 51.867136  | 121.469400 | 2.429913e+10 | 9.986687e+09 | 1.684  |
| min   | 0.000000   | 1.000000   | 1.000000   | 0.000000   | 0.000000e+00 | 0.000000e+00 | 0.000  |
| 25%   | 53.000000  | 107.750000 | 37.750000  | 63.000000  | 3.238344e+09 | 0.000000e+00 | 2.880  |
| 50%   | 106.500000 | 214.500000 | 79.500000  | 143.000000 | 5.801359e+09 | 2.055491e+08 | 5.047  |
| 75%   | 193.250000 | 322.250000 | 132.000000 | 195.000000 | 1.422214e+10 | 9.618685e+08 | 1.143  |
| max   | 300.000000 | 431.000000 | 178.000000 | 771.000000 | 1.810000e+11 | 1.160000e+11 | 1.270  |

```
1 import plotly.express as px # 전체 open close
2 oc = pd.DataFrame(train['OC'].value_counts())
3 oc['status']=oc.index
4 oc.rename(columns={"OC":"count"},inplace=True)
5 fig=px.pie(oc,values="count",names="status",title="Open or Close?",template='seaborn')
6 fig.update_traces(rotation=90,pull=0.05,textinfo="value+percent+label")
7 fig.show()
```

## Open or Close?



```
1 fig,(ax1,ax2) = plt.subplots(1,2,figsize=(20,8))
2 sns.countplot(train['sido'],ax=ax1)
3 ax1.set(title="Distribution of hospitals by region") #지역별 병원 분포
4 ax1.tick_params(labelrotation=60)
5
6 sns.countplot(train['sido'],hue=train['OC'],ax=ax2)
7 plt.xticks(rotation=60)
8 plt.title("Distribution of hospitals by region(+OC)") # 지역별 병원 분포에 open close
```

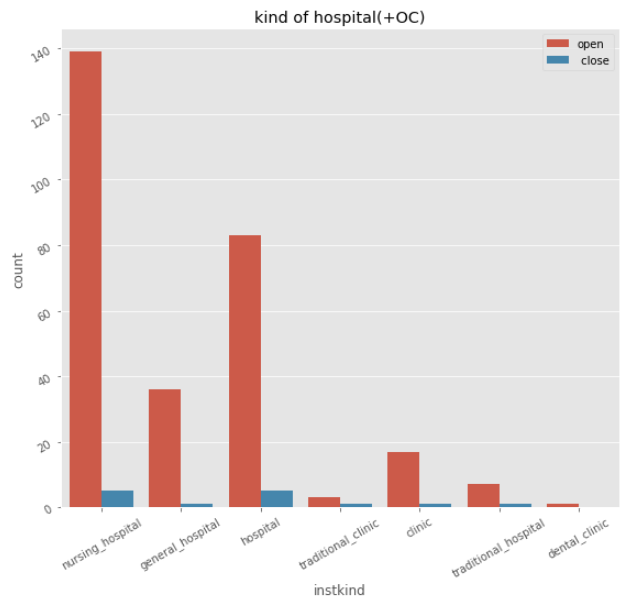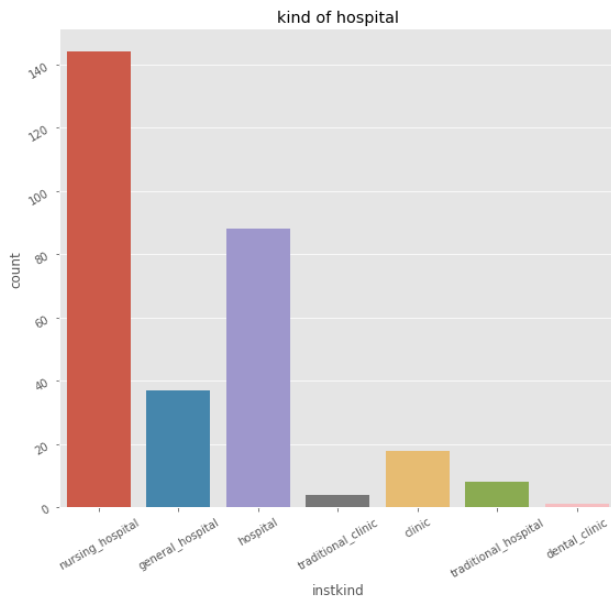Text(0.5, 1.0, 'Distribution of hospitals by region(+OC)')

```
1 fig,(ax1,ax2) = plt.subplots(1,2,figsize=(20,8))
2 sns.countplot(train['instkind'],ax=ax1)
3 ax1.set(title='kind of hospital') # 병원 종류
4 ax1.tick_params(labelrotation=30)
5 sns.countplot(train['instkind'],hue=train['OC'],ax=ax2)
6 ax2.set(title='kind of hospital(+OC)') # 병원 종류 + open close
7 ax2.tick_params(labelrotation=30)
8 plt.legend(loc='upper right')
```

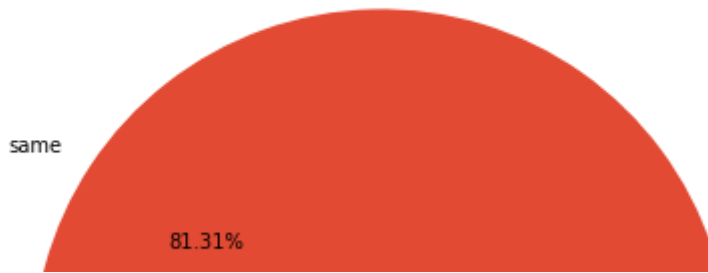<matplotlib.legend.Legend at 0x7fdcf929d6a0>



```
1 plt.figure(figsize=(12,8))
2 data['ownerChange'].value_counts(dropna=False).plot(kind='pie',autopct="%.2f%%")
3 plt.title("Did they change the owner?") # 오너가 바뀌었냐
```
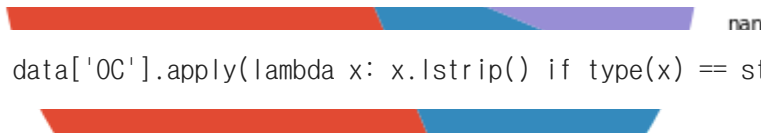
```
Text(0.5, 1.0, 'Did they change the owner?')
```

**Did they change the owner?**

same

81.31%

# PreProcessing

OC ▬▬▬▬▬▬ nan

```
1 data['OC'] = data['OC'].apply(lambda x: x.lstrip() if type(x) == str else x) # close가 안이쁘게
```

```
1 for i in range(data.shape[0]): # employee가 변화 없다고 가정
2     if pd.isnull(data['employee1'][i]) != pd.isnull(data['employee2'][i]): # 직원 하나 있고 하나
3         if pd.isnull(data['employee1'][i]) == True:
4             data['employee1'][i] = data['employee2'][i]
5         elif pd.isnull(data['employee2'][i]) == True:
6             data['employee2'][i] = data['employee1'][i]
7     else:
8         continue
```

```
1 # 직원 수 변화
2 def change_type(x):
3     if type(x) == str:
4         x= x.replace(',','') # ,가 있는게 있음
5         return float(x)
6     else:
7         return x
8 data['employee1'] = data['employee1'].apply(lambda x: change_type(x))
9 data['employee2'] = data['employee2'].apply(lambda x: change_type(x))
10 data['employee_change']=data['employee1']-data['employee2']
```

```
1 import datetime as dt
2 data['openYear'] = data['openDate'].dt.year # 오픈 연도만 저장
3 data.drop('openDate',axis=1,inplace=True) # 자세한 날짜 필요 없으므로 삭제
```

```
1 debt = ['debt1','debt2','liquidLiabilities1','liquidLiabilities2','shortLoan1','shortLoan2','NCl
2       'longLoan1','longLoan2'] # 부채는 모두 마이너스 처리
3 for i in debt:
4     data[i]=data[i].apply(lambda x: 0 if x==0 else -x)
5 data.head()
```

| | index | inst_id | OC | sido | sgg | bedCount | instkind | revenue1 | sa |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 0 | 1 | open | choongnam | 73 | 175.0 | nursing_hospital | 4.217530e+09 | |
| **1** | 1 | 3 | open | gyeongnam | 32 | 410.0 | general_hospital | NaN | |
| **2** | 2 | 4 | open | gyeonggi | 89 | 468.0 | nursing_hospital | 1.004522e+09 | 51 |
| **3** | 3 | 7 | open | incheon | 141 | 353.0 | general_hospital | 7.250734e+10 | |

## ▾ PreProcessing

- Processing Missing value

```
1 data['openYear'].fillna(int(data['openYear'].median()),inplace=True)
```

```
1 data[data['bedCount'].isnull()][['bedCount','instkind']]
```
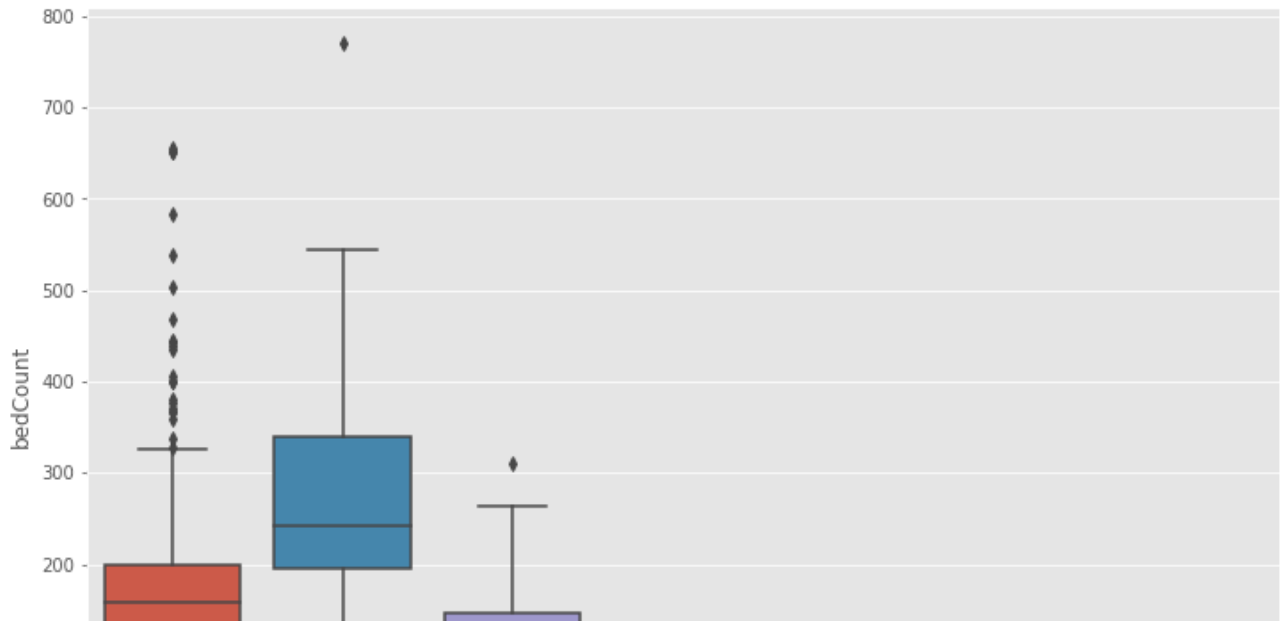
| | bedCount | instkind |
|---|---|---|
| **71** | NaN | traditional_hospital |
| **193** | NaN | NaN |
| **297** | NaN | hospital |
| **298** | NaN | hospital |
| **300** | NaN | traditional_hospital |
| **311** | NaN | traditional_clinic |
| **323** | NaN | hospital |
| **341** | NaN | nursing_hospital |
| **379** | NaN | hospital |
| **385** | NaN | nursing_hospital |
| **400** | NaN | nursing_hospital |
| **424** | NaN | traditional_hospital |
| **426** | NaN | NaN |

```
1 plt.figure(figsize=(12,8))
2 sns.boxplot(x=data['instkind'],y=data['bedCount'])
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fdcfc124320>
```



### #### 이상치가 많으므로 결측치 대체는 중위값으로 한다.

```
1 data['bedCount'].fillna(data.groupby('instkind')['bedCount'].transform('median'),inplace=True)
```

```
1 data[data['bedCount'].isnull()][['bedCount','instkind']] # 둘다 결측치인 것은 영향 안받음
```

```
1 data[data['instkind'].isnull()][['bedCount','instkind']] # instkind 결측치 확인
```

```
1 data.groupby('instkind')['bedCount'].median() # bedCount 가 49는 traditional hospital에 가까움으
```

```
1 data['instkind'].iloc[421] = 'traditional_hospital'
```

```
1 data[['employee1','employee2','bedCount']].corr() # 직원 수와 bedCount의 상관관계가 0.5로 약간
```
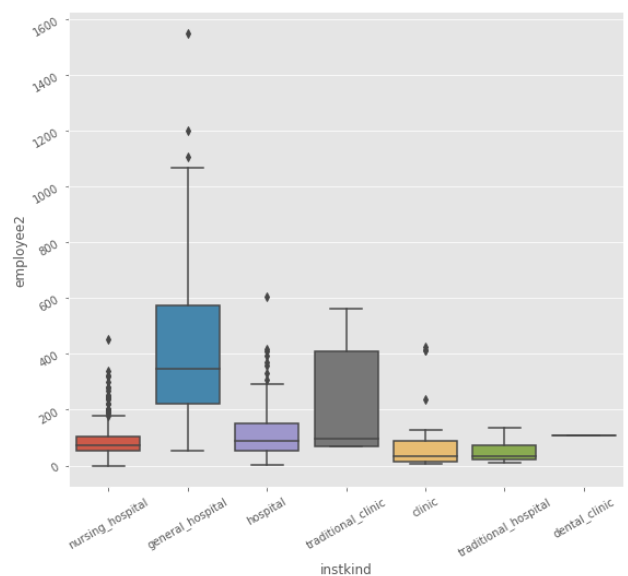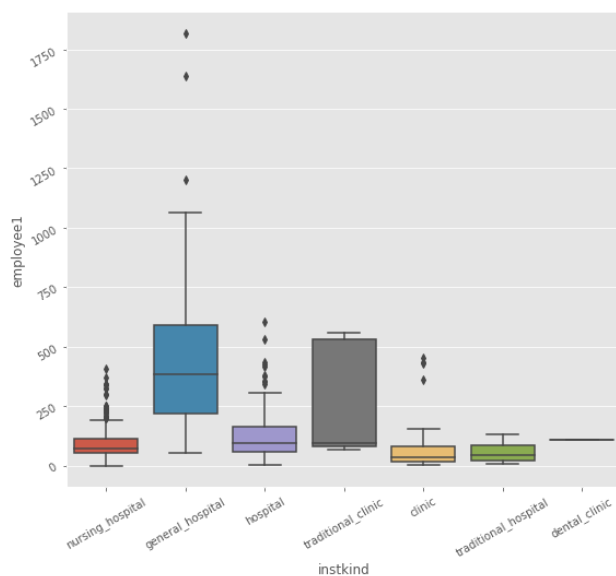
|  | employee1 | employee2 | bedCount |
|---|---|---|---|
| **employee1** | 1.000000 | 0.972793 | 0.505255 |
| **employee2** | 0.972793 | 1.000000 | 0.509500 |
| **bedCount** | 0.505255 | 0.509500 | 1.000000 |

```
1 # 위 상관계수를 기반으로 직원 수에 따라서 bedCount->bedC0unt 와 instkind 상관관계 높음(데이터 숑
2 data[data['instkind'].isnull()]
```

|  | index | inst_id | OC | sido | sgg | bedCount | instkind | revenue1 | salesc |
|---|---|---|---|---|---|---|---|---|---|
| **193** | 193 | 281 | close | gyeonggi | 12 | NaN | NaN | 3.054388e+08 | 2.241614 |
| **426** | 125 | 430 | NaN | jeju | 76 | NaN | NaN | 4.892710e+10 | 4.157148 |

```
1 fig,(ax1,ax2) = plt.subplots(1,2,figsize=(20,8))
2 sns.boxplot(data=data,x='instkind',y='employee1',ax=ax1)
3 sns.boxplot(data=data,x='instkind',y='employee2',ax=ax2)
4 ax1.tick_params(labelrotation=30)
5 ax2.tick_params(labelrotation=30)
```



```
1 data.groupby('instkind')['employee1'].median()
```

```
1 # 15 = clinc
2 # 343 = general_hospital
3 data['instkind'].iloc[193] = 'clinic'
4 data['instkind'].iloc[426] = 'general_hospital'
5 data['bedCount'].iloc[193] = 0 # clinc의 중위값으로 대체
6 data['bedCount'].iloc[426] = 243 # general_hospital 의 중위값으로 대체
```

```
1 data['ownerChange'].fillna(data['ownerChange'].mode()[0],inplace=True) # 최빈값으로 대체
```

- 직원 수 결측치 처리

```
1 df=pd.DataFrame(columns=['col','corr'])
2 for i in range(len(data.corr().keys())):
3     df.loc[i,'col'] = data.corr().keys()[i]
4     df.loc[i,'corr'] = data.corr()['employee1'][i]
5 df
```

| | col | corr |
|---|---|---|
| 0 | index | -0.235603 |
| 1 | inst_id | -0.271752 |
| 2 | sgg | 0.0251274 |
| 3 | bedCount | 0.507025 |
| 4 | revenue1 | 0.88632 |
| 5 | salescost1 | 0.593692 |
| 6 | sga1 | 0.868367 |
| 7 | salary1 | 0.870879 |
| 8 | noi1 | 0.585692 |
| 9 | noe1 | 0.621877 |
| 10 | interest1 | 0.439253 |
| 11 | ctax1 | 0.54814 |
| 12 | profit1 | 0.390353 |
| 13 | liquidAsset1 | 0.743854 |
| 14 | quickAsset1 | 0.740648 |
| 15 | receivableS1 | 0.427205 |
| 16 | inventoryAsset1 | 0.653327 |
| 17 | nonCAsset1 | 0.760217 |
| 18 | tanAsset1 | 0.750643 |
| 19 | OnonCAsset1 | 0.502299 |
| 20 | receivableL1 | -0.0234059 |
| 21 | debt1 | -0.785368 |
| 22 | liquidLiabilities1 | -0.779534 |
| 23 | shortLoan1 | -0.423056 |
| 24 | NCLiabilities1 | -0.648732 |
| 25 | longLoan1 | -0.317576 |
| 26 | netAsset1 | 0.589399 |
| 27 | surplus1 | 0.257669 |
| 28 | revenue2 | 0.879571 |
| 29 | salescost2 | 0.603817 |
| 30 | sga2 | 0.86562 |
| 31 | salary2 | 0.858302 |
| 32 | noi2 | 0.735655 |

| | | |
|---|---|---|
| 32 | noi2 | 0.725655 |
| 33 | noe2 | 0.640054 |
| 34 | interest2 | 0.429575 |
| 35 | ctax2 | 0.460125 |
| 36 | profit2 | 0.424901 |

```
1 corr_features=[]
2 for i in df[abs(df['corr'])>0.7].reset_index()['col']: # 해당년도의 변수만 고려
3     if i[-1] != '2':
4         corr_features.append(i)
5 corr_features
```

```
['revenue1',
 'sga1',
 'salary1',
 'liquidAsset1',
 'quickAsset1',
 'nonCAsset1',
 'tanAsset1',
 'debt1',
 'liquidLiabilities1',
 'employee1']
```

| | | |
|---|---|---|
| 47 | shortLoan2 | -0.3257 |

```
1 emp = data[corr_features]
2 train_emp = emp[emp['employee1'].notnull()]
3 train_emp.dropna(axis=0,inplace=True)
4 train_emp # 회계데이터가 결측치이면 삭제
```

| | revenue1 | sga1 | salary1 | liquidAsset1 | quickAsset1 | nonCAsse |
|---|---|---|---|---|---|---|
| 0 | 4.217530e+09 | 3.961135e+09 | 2.033835e+09 | 1.012700e+09 | 9.976719e+08 | 2.514586e+ |
| 2 | 1.004522e+09 | 4.472197e+08 | 2.964023e+08 | 2.724421e+08 | 2.536822e+08 | 1.204810e+ |
| 3 | 7.250734e+10 | 7.067740e+10 | 3.178605e+10 | 1.304154e+10 | 1.153475e+10 | 4.317936e+ |
| 4 | 4.904354e+10 | 4.765605e+10 | 2.446078e+10 | 6.317084e+09 | 5.873265e+09 | 4.366733e+ |
| 5 | 3.358054e+10 | 2.372791e+10 | 1.665533e+10 | 5.635105e+09 | 5.481680e+09 | 1.864970e+ |
| ... | ... | ... | ... | ... | ... | |
| 421 | 5.583625e+08 | 5.482900e+08 | 2.826852e+08 | 0.000000e+00 | 0.000000e+00 | 0.000000e+ |
| 422 | 4.471030e+08 | 2.581514e+08 | 1.191270e+08 | 2.811359e+08 | 2.336135e+08 | 2.143994e+ |
| 423 | 2.233031e+10 | 1.849255e+10 | 1.232241e+10 | 1.829292e+10 | 1.818429e+10 | 1.307623e+ |
| 424 | 1.833906e+10 | 1.760117e+10 | 6.824241e+09 | 3.706256e+09 | 3.706256e+09 | 7.787147e+ |
| 426 | 4.892710e+10 | 4.721485e+09 | 1.514547e+09 | 1.028286e+10 | 9.002630e+09 | 2.952030e+ |

403 rows × 10 columns

```
1 test_emp = emp[emp['employee1'].isnull()]
2 test_emp
```

```
1 test_emp.drop([test_emp.index[1],test_emp.index[5]],inplace=True)
```

```
1 train_x1 = train_emp.drop(['employee1'],axis=1)
2 train_y1 = train_emp['employee1']
3 test_x1 = test_emp.drop(['employee1'],axis=1)
4 test_y1 = test_emp['employee1']
```

```
1 from sklearn.ensemble import RandomForestRegressor
2 rf=RandomForestRegressor()
3 rf.fit(train_x1,train_y1)
4 pred1 = rf.predict(test_x1)
5 pred1
```

```
array([103.91  ,  81.69 , 435.4    ,  21.91 ,  61.27  , 149.8275,
       115.36  , 228.09 ,  90.33 , 114.19 , 225.2   , 172.652 ,
        44.24 , 140.02 , 124.75  ])
```

```
1 df=pd.DataFrame(columns=['col','corr'])
2 for i in range(len(data.corr().keys())):
3     df.loc[i,'col'] = data.corr().keys()[i]
4     df.loc[i,'corr'] = data.corr()['employee2'][i]
5 df
```

| | col | corr |
|---|---|---|
| **0** | index | -0.252003 |
| **1** | inst_id | -0.284224 |
| **2** | sgg | 0.0232896 |
| **3** | bedCount | 0.511327 |
| **4** | revenue1 | 0.871015 |
| **5** | salescost1 | 0.593576 |
| **6** | sga1 | 0.855941 |
| **7** | salary1 | 0.847373 |
| **8** | noi1 | 0.592255 |
| **9** | noe1 | 0.584647 |
| **10** | interest1 | 0.449565 |
| **11** | ctax1 | 0.487203 |
| **12** | profit1 | 0.325222 |
| **13** | liquidAsset1 | 0.731837 |
| **14** | quickAsset1 | 0.727266 |
| **15** | receivableS1 | 0.470095 |
| **16** | inventoryAsset1 | 0.676041 |
| **17** | nonCAsset1 | 0.785131 |
| **18** | tanAsset1 | 0.778555 |
| **19** | OnonCAsset1 | 0.480342 |

```
1 corr_features2=[]
2 for i in df[abs(df['corr'])>0.7].reset_index()['col']: # 해당년도의 변수만 고려
3     if i[-1] != '1':
4         corr_features2.append(i)
5 corr_features2
```

```
['revenue2',
 'sga2',
 'salary2',
 'noi2',
 'liquidAsset2',
 'quickAsset2',
 'nonCAsset2',
 'tanAsset2',
 'debt2',
 'liquidLiabilities2',
 'employee2']
```

```
1 emp2 = data[corr_features2]
2 train_emp2 = emp2[emp2['employee2'].notnull()]
```

```
3 train_emp2.dropna(axis=0,inplace=True)
4 train_emp2 # 회계데이터가 결측치이면 삭제
```

| | revenue2 | sga2 | salary2 | noi2 | liquidAsset2 | quickAsset |
|---|---|---|---|---|---|---|
| 0 | 4.297848e+09 | 4.057422e+09 | 2.063787e+09 | 16194675.0 | 8.301695e+08 | 8.165705e+0 |
| 2 | 0.000000e+00 | 0.000000e+00 | 0.000000e+00 | 0.0 | 0.000000e+00 | 0.000000e+0 |
| 3 | 6.685834e+10 | 6.492419e+10 | 2.971135e+10 | 476807804.0 | 1.112572e+10 | 9.890540e+0 |
| 4 | 4.808280e+10 | 4.712580e+10 | 2.346004e+10 | 597748128.0 | 4.906776e+09 | 4.464017e+0 |
| 5 | 3.433445e+10 | 2.409622e+10 | 1.638792e+10 | 125681154.0 | 4.869419e+09 | 4.725857e+0 |
| ... | ... | ... | ... | ... | ... | . |
| 421 | 1.160742e+09 | 7.614171e+08 | 4.590994e+08 | 1045466.0 | 0.000000e+00 | 0.000000e+0 |
| 422 | 4.649570e+08 | 3.153264e+08 | 7.356901e+07 | 368343.0 | 3.813002e+07 | 0.000000e+0 |
| 423 | 2.239509e+10 | 1.805503e+10 | 1.138885e+10 | 483447584.0 | 1.708819e+10 | 1.701879e+1 |
| 424 | 1.911503e+10 | 1.824717e+10 | 7.250614e+09 | 41234195.0 | 4.588941e+09 | 4.588941e+0 |
| 426 | 4.758477e+10 | 5.061219e+09 | 1.404341e+09 | 99795507.0 | 7.848233e+09 | 6.707102e+0 |

403 rows × 11 columns

```
1 test_emp2 = emp2[emp2['employee2'].isnull()]
2 test_emp2.drop([test_emp2.index[1],test_emp2.index[5]],inplace=True)
```

```
1 train_x2 = train_emp2.drop(['employee2'],axis=1)
2 train_y2 = train_emp2['employee2']
3 test_x2 = test_emp2.drop(['employee2'],axis=1)
4 test_y2 = test_emp2['employee2']
```

```
1 rf.fit(train_x2,train_y2)
2 pred2 = rf.predict(test_x2)
3 pred2
```

```
1 idx = list(data[data['employee1'].isnull()].index)
2 idx.remove(60)
3 idx.remove(258)
4 idx
```

```
1 data['employee1'][idx] = pred1
2 data['employee2'][idx] = pred2
```

```
1 data.groupby('instkind')['employee1','employee2','bedCount'].median() # 두개 다 결측치인 것은 병
```

| instkind | employee1 | employee2 | bedCount |
|---|---|---|---|
| clinic | 28.0 | 27.5 | 0.0 |
| dental_clinic | 107.0 | 109.0 | 0.0 |
| general_hospital | 377.0 | 346.0 | 243.0 |
| hospital | 96.5 | 90.0 | 85.0 |
| nursing_hospital | 74.5 | 74.0 | 157.5 |
| traditional_clinic | 86.0 | 82.5 | 0.0 |
| traditional_hospital | 42.0 | 34.0 | 44.0 |

```
1 data['employee1'].iloc[60] = 75
2 data['employee2'].iloc[60] = 74
3 data['employee1'].iloc[258] = 28
4 data['employee2'].iloc[258] = 28
```

```
1 data['employee1'] = data['employee1'].apply(lambda x: round(x)) # 직원 수 반올림
2 data['employee2'] = data['employee2'].apply(lambda x: round(x))
```

```
1 data['employee_change'] = data['employee1'] - data['employee2'] #  update employee_change columr
```

## Preprocessing

- 회계데이터 결측치 처리

```
 1 tmp2016=[]
 2 tmp2017=[]
 3 for i in data.columns:
 4   if i[-1] == '2':
 5     tmp2016.append(i)
 6   elif i[-1] == '1':
 7     tmp2017.append(i)
 8 tmp2016.remove('employee2')
 9 tmp2017.remove('employee1')
10 acc2016 = data[tmp2016]
11 acc2017 = data[tmp2017]
```

```
1 zero_idx2016=[]
2 zero_idx2017=[]
3 for i in range(acc2016.shape[0]):
4   if sum(acc2016.iloc[i].values==0) == 0:
5     zero_idx2016.append(i)
6 for j in range(acc2017.shape[0]):
```

```
 7    if sum(acc2017.iloc[j].values==0) == 0:
 8      zero_idx2017.append(j)
 9 print('2016: ',zero_idx2016) # 2016년 회계 데이터 전체가 0인 index들
10 print('2017: ',zero_idx2017) # 2017년 회계 데이터 전체가 0인 index들
11 # 회계데이터가 모두 0인 데이터는 없습니다. 나온 결과값 모두 결측치 값입니다.
```

```
2016:  [1, 14, 32, 55, 60, 248, 257, 258, 316, 357]
2017:  [1, 14, 32, 55, 60, 248, 257, 258, 316, 357]
```

```
1 data['old'] = 2018-data['openYear'] # 2018년 데이터이므로 오픈한지 얼마나 지났는지
```

```
1 from xgboost import XGBRegressor # xgb로 회계데이터 모델링
2 xgb = XGBRegressor(n_estimators=1000)
```

```
1 df= data.drop(['index','inst_id','OC','sido','sgg','ownerChange','openYear','instkind','employee
```

```
1 idx=list(data.index)
2 for i in zero_idx2016:
3   idx.remove(i)
```

```
1 na_col = []
2 for i,j in zip(tmp2016,tmp2017):
3   na_col.append(i)
4   na_col.append(j)
5 na_col # 결측치가 있는 데이터 모아놓기
```

```
['revenue2',
 'revenue1',
 'salescost2',
 'salescost1',
 'sga2',
 'sga1',
 'salary2',
 'salary1',
 'noi2',
 'noi1',
 'noe2',
 'noe1',
 'interest2',
 'interest1',
 'ctax2',
 'ctax1',
 'profit2',
 'profit1',
 'liquidAsset2',
 'liquidAsset1',
 'quickAsset2',
```

```
1 for i in na_col:
2   train_x = df.drop(i,axis=1).iloc[idx]
3   train_y = df[i].iloc[idx]
4   test_x = df.drop(i,axis=1).iloc[zero_idx2016]
5   xgb.fit(train_x,train_y)
6   pred = xgb.predict(test_x)
7   data[i][zero_idx2016] = pred
8 data.isnull().sum()
```

```
[05:57:11] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:13] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:14] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:15] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:16] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:17] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:19] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:20] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:21] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:22] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:23] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:24] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:26] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:27] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:28] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:29] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:30] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:32] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:33] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:34] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:35] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:36] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:38] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:39] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:40] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:41] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:42] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:43] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:45] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:46] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:47] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:48] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
[05:57:49] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
```

# Preprocessing

- final

```
[05:57:58] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
```

```
1 data['profit_diff'] = data['profit1'] − data['profit2']
2 data['debt_diff'] = data['debt1'] − data['debt2']
3 data['surplus_diff'] = data['surplus1'] − data['surplus2']
4 data['netasset_diff'] = data['netAsset1'] − data['netAsset2']
```

```
[05:58:05] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now depreca
```

```
1 dic={'open':1,'close':0}
2 data['OC'] = data['OC'].map(dic)
```

```
1 dic={'same':1,'change':0}
2 data['ownerChange'] = data['ownerChange'].map(dic)
```

```
1 data.drop(['sgg','openYear','index'],axis=1,inplace=True) # 필요없는 것 다른 것으로 대체할 수 있
```

```
1 # 충남 충북 -> 충청
```

```
2 # 경남 경북 -> 경상
3 # 전남 전북 -> 전라
4 # 경기 -> 경기
5 # 강원 -> 강원
6 # 제주 -> 제주
7 # 인천 부산 대구 울산 세종 대전 광주 ->광역시
```

```
1 sido_dic={"choongnam":"충청","choongbuk":"충청","gyeongnam":"경상","gyeongbuk":"경상", "jeonnam'
2          'seoul':'서울','gangwon':'강원','jeju':'제주',"sejong":"광역시","daejeon":"광역시",
3          "busan":"광역시","daegu":"광역시","ulsan":"광역시","gwangju":'광역시','incheon':"광역시
4 data['sido'] = data.sido.map(sido_dic)
```

```
1 from sklearn.preprocessing import OneHotEncoder # onehotencoding dataframe 변환 함수
2 def ohe_trans(data,col):
3     ohe=OneHotEncoder()
4     x= ohe.fit_transform(data[col].values.reshape(-1,1)).toarray()
5     tp = []
6     for i in range(data[col].unique().size): # onehot 컬럼 생성
7         tp.append(col[0]+str(i))
8     ohe_df = pd.DataFrame(x,columns = tp)
9     return ohe_df
10
```

```
1 ohe_sido = ohe_trans(data,'sido')
2 ohe_kind = ohe_trans(data,'instkind')
3 data.drop(['sido','instkind'],axis=1,inplace=True)
```

```
1 data=pd.concat([data,ohe_sido,ohe_kind],axis=1)
```

```
1 train_data = data.iloc[:train_len]
2 test_data = data.iloc[train_len:]
```

```
1 train_data.drop('inst_id',axis=1,inplace=True)
2 test_data.drop('inst_id',axis=1,inplace=True)
```

## MODELING

- use ML(XGBCLASSIFIER,RANDOMFORESTCLASSIFIER)

```
1 from sklearn.ensemble import RandomForestClassifier
2 from xgboost import XGBClassifier
3 rf = RandomForestClassifier(n_estimators = 1000)
4 xgb = XGBClassifier(n_estimators = 1000)
```

```
1 train_x = train_data.drop('OC',axis=1)
2 train_y = train_data['OC']
3 test_x = test_data.drop('OC',axis=1)
4 test_y = test_data['OC']
```

4 test_y = test_data['OC']

```
1 rf.fit(train_x,train_y)
2 pred_rf = rf.predict(test_x)
3 xgb.fit(train_x,train_y)
4 pred_xgb = xgb.predict(test_x)
```
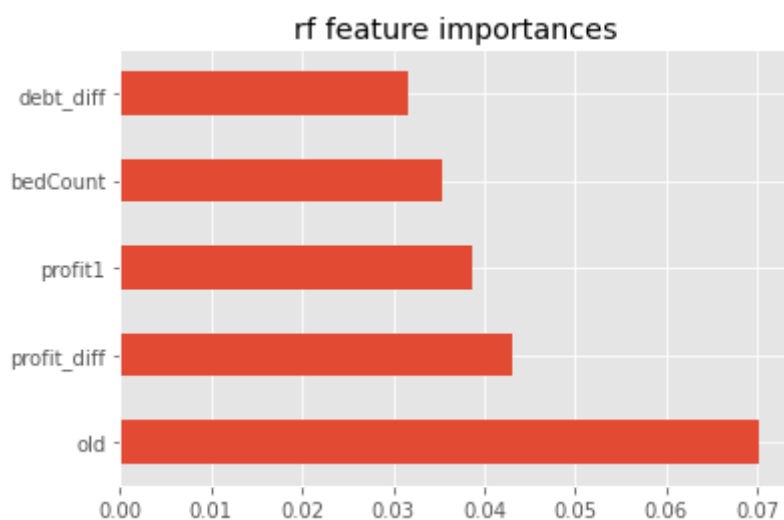
```
1 answer['OC'] = pred_rf
2 #answer.to_csv('rf_predict(hos oc).csv',index=False)
```

```
1 answer['OC'] = pred_xgb
2 #answer.to_csv('xgb_predict(hos oc).csv',index=False)
```

| 432657 | xgb score(hos_oc).csv | 2020-06-07 13:28:54 | 0.8503937008 |
|--------|------------------------|---------------------|--------------|
| 432654 | rf score(hos_oc).csv | 2020-06-07 13:07:02 | 0.842519685 |

```
1 feat_importance=pd.Series(rf.feature_importances_,index=train_x.columns)
2 feat_importance.nlargest(5).plot(kind='barh')
3 plt.title('rf feature importances')
4
```
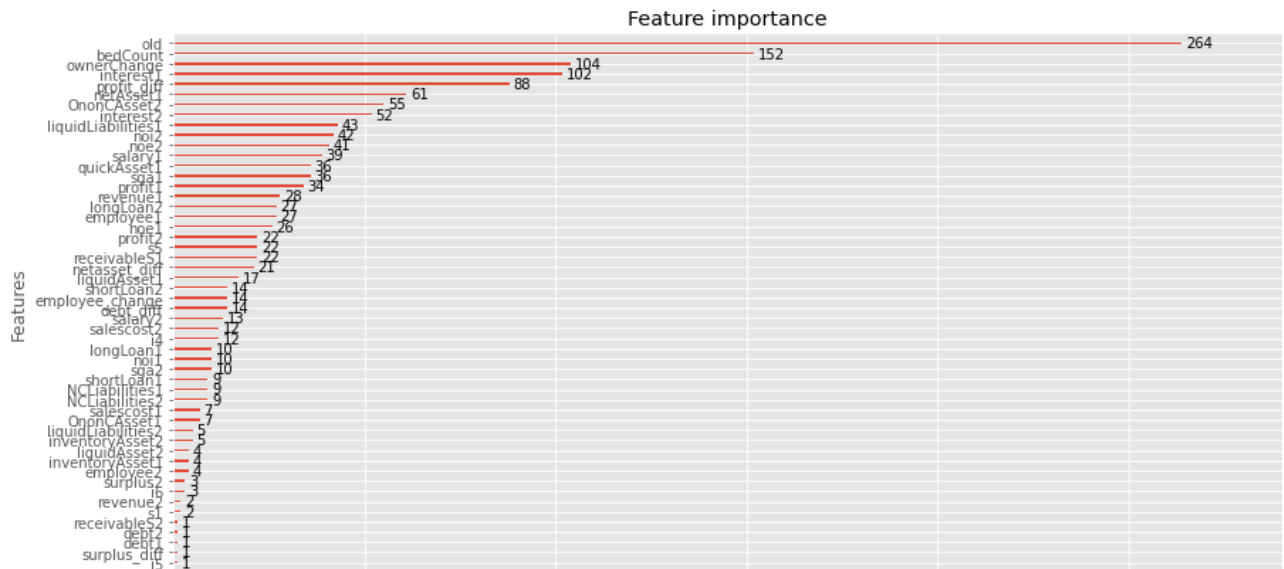
Text(0.5, 1.0, 'rf feature importances')



```
1 from xgboost import plot_importance
2 plt.rcParams["figure.figsize"] = (14, 7)
3 plot_importance(xgb)
```

```
<matplotlib.axes._subplots.AxesSubplot at 0x7fdc43231b70>
```


Feature importance

- 두 모델 공통적으로 old 변수가 OC를 가장 판단하는데 가장 중요한 변수로 작용하였다. 이 어서 profit_diff,bedCount 등의 변수가 예측하는데 중요한 변수로 되었다.

# ▾ MODELING

- 인공신경망

```
 1 import torch
 2 import torch.nn as nn
 3 import torch.optim as optim
 4 import torch.nn.init as init
 5
 6 # + CNN할때 필요한 것
 7 import torch.utils.data as data
 8 import torchvision.datasets as dset
 9 import torchvision.transforms as transforms
10 from torch.utils.data import DataLoader
```

```
 1 device = torch.device("cuda:0" if torch.cuda.is_available() else "cpu")
```

```
 1 train_df = train_data.drop(['OC'],axis=1)
 2 train_OC = train_data['OC']
```

```
 1 from sklearn.model_selection import train_test_split # 검정용 데이터 생성
 2 valid_train_x,valid_x,valid_train_y,valid_y = train_test_split(train_df,train_OC,test_size=.2,ra
 3
 4 valid_train_tensor = torch.tensor(valid_train_x.values,dtype=torch.float) # 검정용 데이터 train
 5 valid_train_tensor_y = torch.tensor(valid_train_y.values,dtype=torch.float) # 검정용 데이터 trai
 6
 7 train_tensor = valid_train_tensor.type(torch.FloatTensor).to(device)
 8 train_tensor_y = valid_train_tensor_y.type(torch.LongTensor).to(device)
 9
10 valid_test_tensor = torch.tensor(valid_x.values,dtype=torch.float).to(device)
```

```
1 num_epoch = 10000
2 lr = [0.005,0.002,0.0005,0.0002]
3 s_lr = [0.05,0.02,0.005,0.002] # 시그모이드는 학습률 크게
```

```
1 # 활성화 함수를 변경 시키면서 성능 확인 필요
2
3 w = 73  # input의 개수
4 model = nn.Sequential(
5         nn.Linear(1*w,6*w),
6         nn.ReLU(),
7         nn.Linear(6*w,10*w),
8         nn.ReLU(),
9         nn.Linear(10*w,6*w),
10         nn.ReLU(),
11         nn.Linear(6*w,1*w),
12         nn.LeakyReLU(),
13         nn.Linear(w,2),
14     )
15
16 loss_func = nn.CrossEntropyLoss()
17 # optimizer = optim.Adam(model.parameters(),lr=0.002)
18 #optimizer2 = optim.Adam(model.parmeters(),lr=0.0002)
19 model.to(device)
```

```
1 # 모델 함수로 만들어서 계속 초기화 할 수 있도록
2
3 def model_RRRLR():
4     w = 73  # input의 개수
5     model = nn.Sequential(
6         nn.Linear(1*w,6*w), nn.ReLU(),
7         nn.Linear(6*w,10*w), nn.ReLU(),
8         nn.Linear(10*w,6*w), nn.ReLU(),
9         nn.Linear(6*w,1*w), nn.LeakyReLU(),
10         nn.Linear(w,2) )
11     loss_func = nn.CrossEntropyLoss()
12     return model.to(device)
13
14 def model_RRRR():
15     w = 73  # input의 개수
16     model = nn.Sequential(
17         nn.Linear(1*w,6*w), nn.ReLU(),
```

```
18          nn.Linear(6*w,10*w), nn.ReLU(),
19          nn.Linear(10*w,6*w), nn.ReLU(),
20          nn.Linear(6*w,1*w), nn.ReLU(),
21          nn.Linear(w,2) )
22      loss_func = nn.CrossEntropyLoss()
23      return model.to(device)
24
25 def model_RRRT():
26      w = 73  # input의 개수
27      model = nn.Sequential(
28          nn.Linear(1*w,6*w), nn.ReLU(),
29          nn.Linear(6*w,10*w), nn.ReLU(),
30          nn.Linear(10*w,6*w), nn.ReLU(),
31          nn.Linear(6*w,1*w), nn.Tanh(),
32          nn.Linear(w,2) )
33      loss_func = nn.CrossEntropyLoss()
34      return model.to(device)
35
36 def model_SSSS():
37      w = 73  # input의 개수
38      model = nn.Sequential(
39          nn.Linear(1*w,6*w), nn.Sigmoid(),
40          nn.Linear(6*w,10*w), nn.Sigmoid(),
41          nn.Linear(10*w,6*w), nn.Sigmoid(),
42          nn.Linear(6*w,1*w), nn.Sigmoid(),
43          nn.Linear(w,2) )
44      loss_func = nn.CrossEntropyLoss()
45      return model.to(device)
46
47 def model_TTTT():
48      w = 73  # input의 개수
49      model = nn.Sequential(
50          nn.Linear(1*w,6*w), nn.Tanh(),
51          nn.Linear(6*w,10*w), nn.Tanh(),
52          nn.Linear(10*w,6*w), nn.Tanh(),
53          nn.Linear(6*w,1*w), nn.Tanh(),
54          nn.Linear(w,2) )
55      loss_func = nn.CrossEntropyLoss()
56      return model.to(device)
57
58 def model_LRLRLRLR():
59      w = 73  # input의 개수
60      model = nn.Sequential(
61          nn.Linear(1*w,6*w), nn.LeakyReLU(),
62          nn.Linear(6*w,10*w), nn.LeakyReLU(),
63          nn.Linear(10*w,6*w), nn.LeakyReLU(),
64          nn.Linear(6*w,1*w), nn.LeakyReLU(),
65          nn.Linear(w,2) )
66      loss_func = nn.CrossEntropyLoss()
67      return model.to(device)
```

# ▾ ReLU - ReLU - RELU - LeakyReLU

```python
1 loss_array_array = []
2 model_array = []
3 for r in lr:
4     loss_array = []
5     model = model_RRRLR()
6     optimizer = optim.Adam(model.parameters(),r)
7     print("학습률 : {}".format(r))
8     for i in range(num_epoch):
9         optimizer.zero_grad()
10         output = model(train_tensor).to(device)
11         loss = loss_func(output,train_tensor_y)
12         if i%1000==0:
13             print(loss)
14         loss.backward()
15         optimizer.step()
16         loss_array.append(loss)
17     model_array.append(model)
18     loss_array_array.append(loss_array)
```
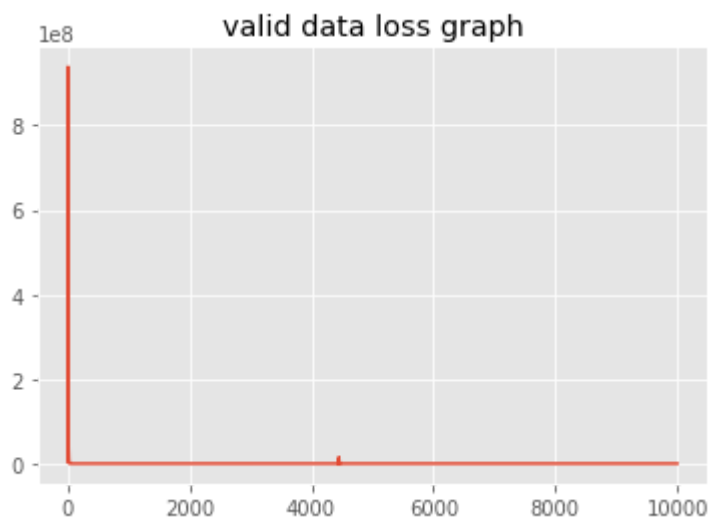
학습률 : 0.005
tensor(1828784.7500, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(96.9898, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(37.1877, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(147.8693, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(88.9446, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(161.2969, device='cuda:0', grad_fn=<NllLossBackward>)

```
1 for loss_array in loss_array_array:
2    plt.plot(loss_array)
3    plt.title("valid data loss graph")
4    plt.show()
```

valid data loss graph

```
1 output_array = []
2 for m in model_array:
3     output = m(valid_test_tensor)
4     output_array.append(output)
5 # output
```

```
1 valid_result_array=[]
2 for output in output_array:
3     valid_result = []
4     for i in output:
5         if i.argmax() == 1:
6             valid_result.append(1)
7         else:
8             valid_result.append(0)
9     valid_result_array.append(valid_result)
```

```
1 from sklearn.metrics import accuracy_score # 정확도 측정
2 for valid_result in valid_result_array:
3     print(accuracy_score(valid_y,valid_result)) # 오류나면 데이터type 통일해주기
```

## ▾ ReLU - ReLU - ReLU - ReLU

```
1 loss_array_array = []
2 model_array = []
3 for r in lr:
4     loss_array = []
5     model = model_RRRR()
6     optimizer = optim.Adam(model.parameters(),r)
7     print("학습률 : {}".format(r))
8     for i in range(num_epoch):
9         optimizer.zero_grad()
10        output = model(train_tensor).to(device)
```
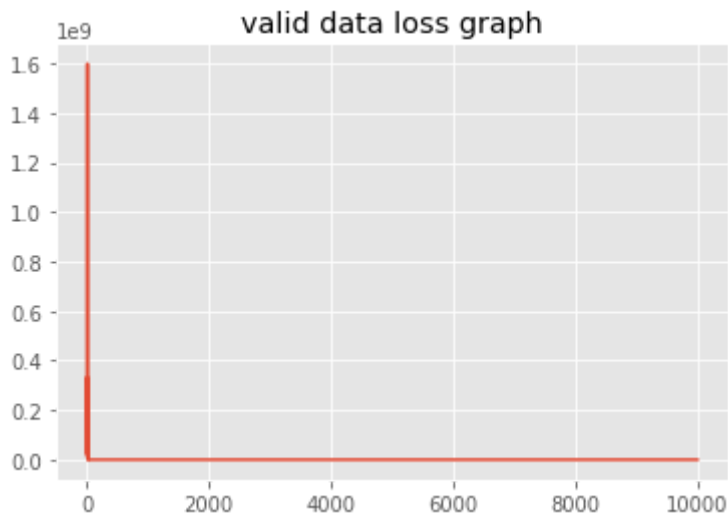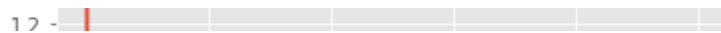
```
10      output = model(train_tensor).to(device)
11      loss = loss_func(output,train_tensor_y)
12      if i%5000==0:
13          print(loss)
14      loss.backward()
15      optimizer.step()
16      loss_array.append(loss)
17   model_array.append(model)
18   loss_array_array.append(loss_array)
```

```
1 for loss_array in loss_array_array:
2    plt.plot(loss_array)
3    plt.title("valid data loss graph")
4    plt.show()
```

valid data loss graph



valid data loss graph

valid data loss graph

```
1 output_array = []
2 for m in model_array:
3     output = m(valid_test_tensor)
4     output_array.append(output)
```

```
1 valid_result_array=[]
2 for output in output_array:
3     valid_result = []
4     for i in output:
5         if i.argmax() == 1:
6             valid_result.append(1)
7         else:
8             valid_result.append(0)
9     valid_result_array.append(valid_result)
```
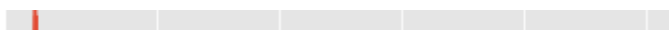
```
1 for valid_result in valid_result_array:
2     print(accuracy_score(valid_y,valid_result)) # 오류나면 데이터type 통일해주기
```

0.9672131147540983
0.9508196721311475
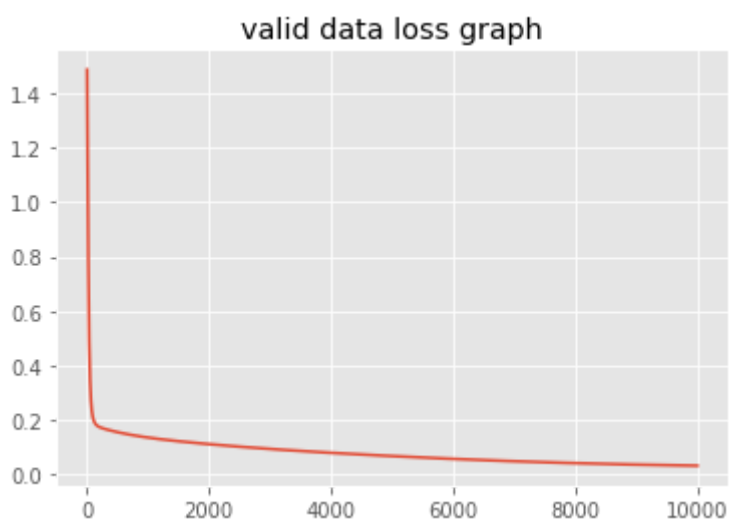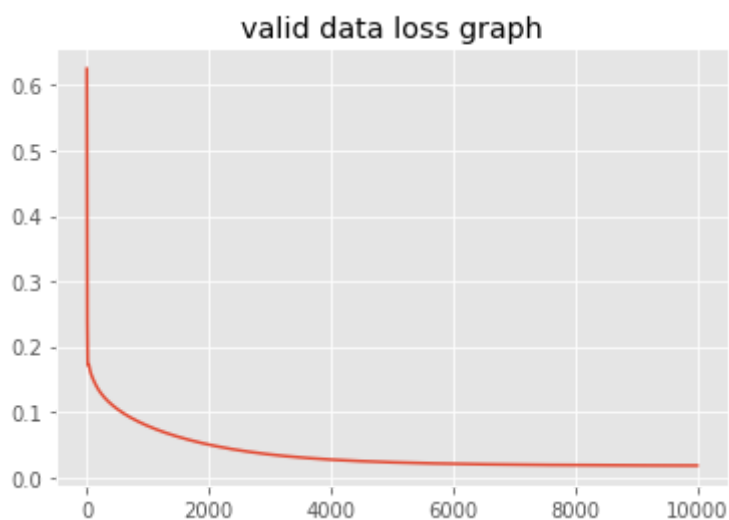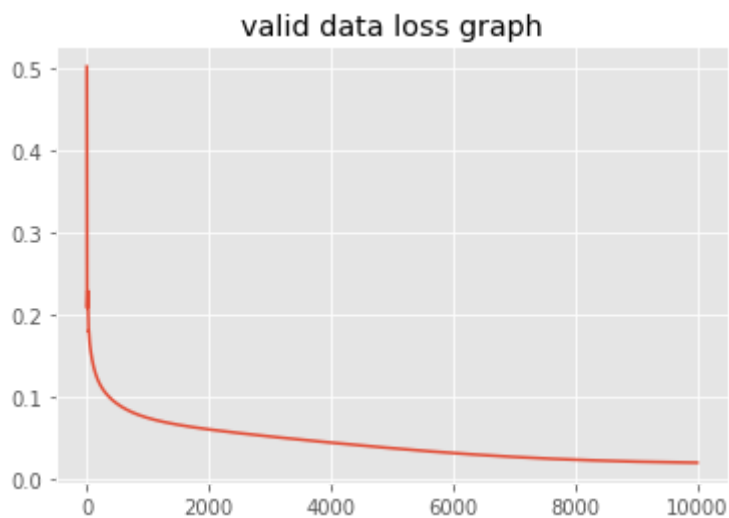0.9672131147540983
0.9672131147540983

## ▾ ReLU - ReLU - ReLU - Tanh

```
1 loss_array_array = []
2 model_array = []
3 for r in lr:
4     loss_array = []
5     model = model_RRRT()
6     optimizer = optim.Adam(model.parameters(),r)
7     print("학습률 : {}".format(r))
8     for i in range(num_epoch):
9         optimizer.zero_grad()
10        output = model(train_tensor).to(device)
11        loss = loss_func(output,train_tensor_y)
12        if i%5000==0:
13            print(loss)
14        loss.backward()
15        optimizer.step()
16        loss_array.append(loss)
17    model_array.append(model)
18    loss_array_array.append(loss_array)
```

```
1 for loss_array in loss_array_array:
2     plt.plot(loss_array)
3     plt.title("valid data loss graph")
4     plt.show()
```

## valid data loss graph



## valid data loss graph



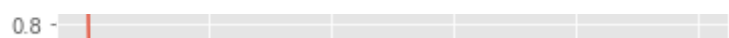## valid data loss graph



```
1 output_array = []
2 for m in model_array:
3     output = m(valid_test_tensor)
4     output_array.append(output)
```



```
1 valid_result_array=[]
2 for output in output_array:
3     valid_result = []
4     for i in output:
5         if i.argmax() == 1:
6             valid_result.append(1)
7         else:
```

```
8            valid_result.append(0)
9        valid_result_array.append(valid_result)
```

```
1 for valid_result in valid_result_array:
2     print(accuracy_score(valid_y,valid_result)) # 오류나면 데이터type 통일해주기
```

## ▼ Sigmoid - Sigmoid - Sigmoid - Sigmoid

```
1 loss_array_array = []
2 model_array = []
3 for r in s_lr:
4     loss_array = []
5     model = model_RRRT()
6     optimizer = optim.Adam(model.parameters(),r)
7     print("학습률 : {}".format(r))
8     for i in range(num_epoch):
9         optimizer.zero_grad()
10        output = model(train_tensor).to(device)
11        loss = loss_func(output,train_tensor_y)
12        if i%5000==0:
13            print(loss)
14        loss.backward()
15        optimizer.step()
16        loss_array.append(loss)
17    model_array.append(model)
18    loss_array_array.append(loss_array)
```

```
1 for loss_array in loss_array_array:
2     plt.plot(loss_array)
3     plt.title("valid data loss graph")
4     plt.show()
```

```
1 output_array = []
2 for m in model_array:
3     output = m(valid_test_tensor)
4     output_array.append(output)
```

```
1 valid_result_array=[]
2 for output in output_array:
3     valid_result = []
4     for i in output:
5         if i.argmax() == 1:
6             valid_result.append(1)
7         else:
8             valid_result.append(0)
9     valid_result_array.append(valid_result)
```

```
1 for valid_result in valid_result_array:
2     print(accuracy_score(valid_y,valid_result)) # 오류나면 데이터type 통일해주기
```

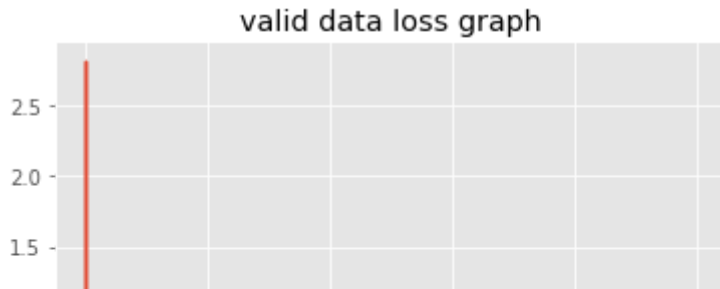## ▾ Tanh - Tanh - Tanh - Tanh

```
1 loss_array_array = []
2 model_array = []
3 for r in lr:
4     loss_array = []
5     model = model_TTTT()
6     optimizer = optim.Adam(model.parameters(),r)
7     print("학습률 : {}".format(r))
8     for i in range(num_epoch):
9         optimizer.zero_grad()
10        output = model(train_tensor).to(device)
11        loss = loss_func(output,train_tensor_y)
12        if i%5000==0:
13            print(loss)
14        loss.backward()
15        optimizer.step()
16        loss_array.append(loss)
17    model_array.append(model)
18    loss_array_array.append(loss_array)
```

```
1 for loss_array in loss_array_array:
2     plt.plot(loss_array)
3     plt.title("valid data loss graph")
4     plt.show()
```

## valid data loss graph



```
1 output_array = []
2 for m in model_array:
3     output = m(valid_test_tensor)
4     output_array.append(output)
```



```
1 valid_result_array=[]
2 for output in output_array:
3     valid_result = []
4     for i in output:
5         if i.argmax() == 1:
6             valid_result.append(1)
7         else:
8             valid_result.append(0)
9     valid_result_array.append(valid_result)
```



```
1 for valid_result in valid_result_array:
2     print(accuracy_score(valid_y,valid_result)) # 오류나면 데이터type 통일해주기
```

```
0.9180327868852459
0.9508196721311475
0.9344262295081968
0.9508196721311475
```



# ▼ LeakyReLU - LeakyReLU - LeakyReLU - LeakyReLU



```
1 loss_array_array = []
2 model_array = []
3 for r in lr:
4     loss_array = []
5     model = model_LRLRLRLR()
6     optimizer = optim.Adam(model.parameters(),r)
7     print("학습률 : {}".format(r))
8     for i in range(num_epoch):
9         optimizer.zero_grad()
10        output = model(train_tensor).to(device)
11        loss = loss_func(output,train_tensor_y)
12        if i%5000==0:
13            print(loss)
14        loss.backward()
15        optimizer.step()
16        loss_array.append(loss)
17    model_array.append(model)
18    loss array array.append(loss array)
```

학습률 : 0.005
tensor(2566436.2500, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(267.5387, device='cuda:0', grad_fn=<NllLossBackward>)
학습률 : 0.002
tensor(11951636., device='cuda:0', grad_fn=<NllLossBackward>)
tensor(902.5310, device='cuda:0', grad_fn=<NllLossBackward>)
학습률 : 0.0005
tensor(14195371., device='cuda:0', grad_fn=<NllLossBackward>)
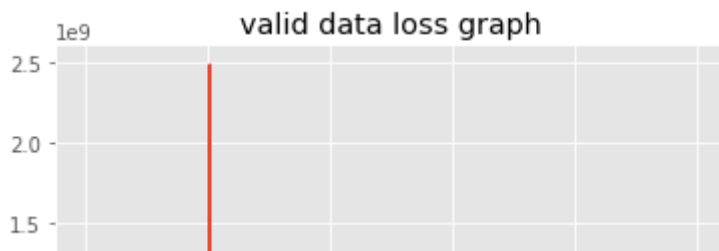tensor(1443517.8750, device='cuda:0', grad_fn=<NllLossBackward>)
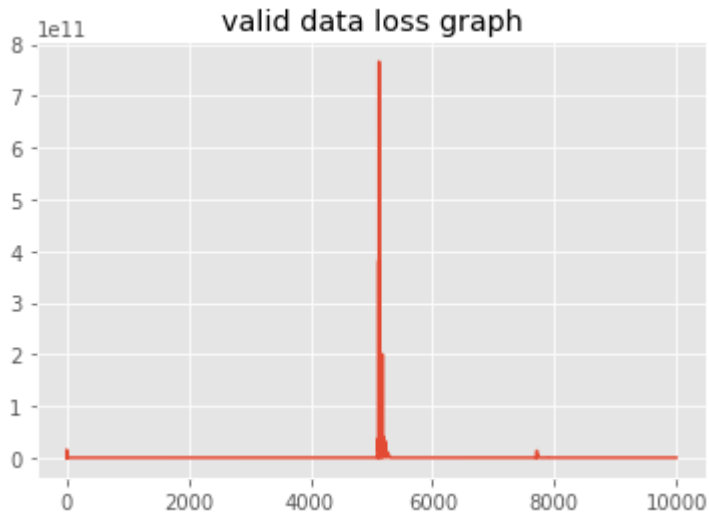학습률 : 0.0002
tensor(2753409.5000, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(2489.4126, device='cuda:0', grad_fn=<NllLossBackward>)

```
1 for loss_array in loss_array_array:
2     plt.plot(loss_array)
3     plt.title("valid data loss graph")
4     plt.show()
```

valid data loss graph



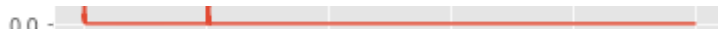valid data loss graph

```
1 output_array = []
2 for m in model_array:
3     output = m(valid_test_tensor)
4     output_array.append(output)
```
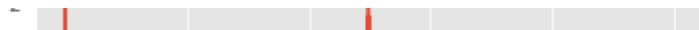


```
1 valid_result_array=[]
2 for output in output_array:
3     valid_result = []
4     for i in output:
5         if i.argmax() == 1:
6             valid_result.append(1)
7         else:
8             valid_result.append(0)
9     valid_result_array.append(valid_result)
```
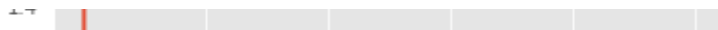


```
1 for valid_result in valid_result_array:
2     print(accuracy_score(valid_y,valid_result)) # 오류나면 데이터type 통일해주기
```

0.9344262295081968
0.9508196721311475
0.9672131147540983
0.9672131147540983



# 정확도 측정 결과

| | RRRLR | RRRR | RRRT | TTTT | LRLRLRLR | | SSSS |
|---|---|---|---|---|---|---|---|
| 0.005 | 0.93 | 0.95 | 0.86 | 0.967 | 0.93 | 0.05 | 0.91 |

# Modeling test data

| 0.0002 | 0.907 | 0.95 | 0.808 | 0.95 | 0.907 | 0.002 | 0.91 |

```
1
2 train_df_tensor = torch.tensor(train_df.values,dtype=torch.float)
3 train_data_Output = torch.tensor(train_OC.values, dtype=torch.float)
4 test_df = test_data.drop('OC',axis=1)
5 test_df_tensor = torch.tensor(test_df.values,dtype=torch.float)
```

```
1 num_epoch = 10000
2
3 # train_tensor = train_tensor.type(torch.FloatTensor)
4 #train_tensor = train_tensor.cuda() # GPU로 보냄
5 #noise = init.normal_(train_tensor,std=1)
6 # x = init.uniform_(train_tensor,-15,15)
7 lr = [0.005,0.002,0.0005,0.0002]
```

```
 1 # # 아래 코드는 특성의 개수가 73 -> 6*73 -> 10*73 -> 6*73 -> 73 -> 1개로 변하는 인공신경망입니디
 2
 3 # w = 73   # input의 개수
 4 # model = nn.Sequential(
 5 #           nn.Linear(1*w,6*w),
 6 #           nn.ReLU(),
 7 #           nn.Linear(6*w,10*w),
 8 #           nn.ReLU(),
 9 #           nn.Linear(10*w,6*w),
10 #           nn.ReLU(),
11 #           nn.Linear(6*w,1*w),
12 #           nn.LeakyReLU(),
13 #           nn.Linear(w,2),
14 #       )
15
16 # loss_func = nn.CrossEntropyLoss()
17 # optimizer = optim.Adam(model.parameters(),lr=0.002)
18 # #optimizer2 = optim.Adam(model.parmeters(),lr=0.0002)
19 # model.to(device)
```

```
1 final_train_tensor = torch.tensor(train_df.values,dtype=torch.float).to(device)
2 train_y_tensor = torch.tensor(train_OC.values,dtype=torch.float).to(device)
3 final_train_tensor = final_train_tensor.type(torch.FloatTensor).to(device)
4 train_y_tensor = train_y_tensor.type(torch.LongTensor).to(device)
```

```
1 test_x = test_data.drop('OC',axis=1)
2 test_tensor = torch.tensor(test_x.values,dtype=torch.float).to(device)
```

```
1 model = model_RRRLR()
2 loss_func = nn.CrossEntropyLoss()
```

```
 3 optimizer = optim.Adam(model.parameters(),lr=0.0005)
 4 model.to(device)
 5 loss_array = []
 6 for i in range(num_epoch):
 7     optimizer.zero_grad()
 8     output = model(final_train_tensor).to(device)
 9     loss = loss_func(output,train_y_tensor)
10     if i%1000==0:
11         print(loss)
12     loss.backward()
13     optimizer.step()
14     loss_array.append(loss)
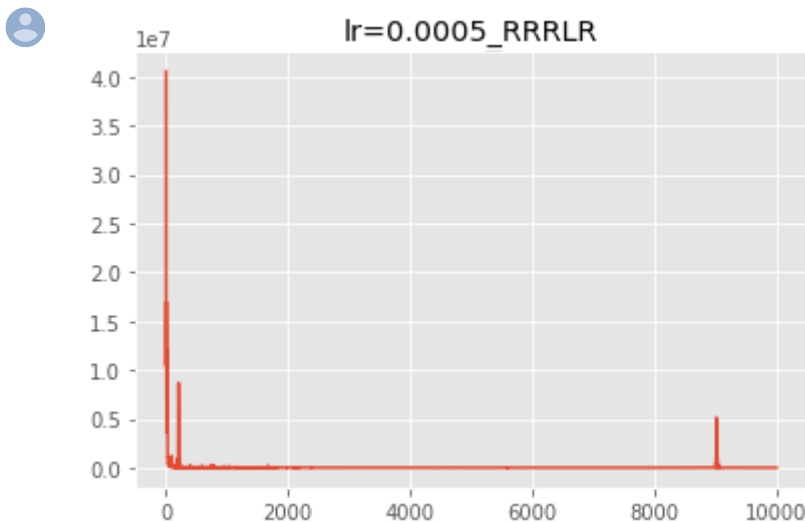```

```
tensor(40525448., device='cuda:0', grad_fn=<NllLossBackward>)
tensor(5737.2456, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(19227.3457, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(8675.5332, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(2017.1469, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(958.8677, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(486.0355, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(70.5464, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(2.4756, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(14675.6328, device='cuda:0', grad_fn=<NllLossBackward>)
```

```
 1 # 손실 그래프
 2 plt.plot(loss_array)
 3 plt.title("lr=0.0005_RRRLR")
 4 plt.show()
 5 print(loss_array[-1])
```



```
tensor(218.8876, device='cuda:0', grad_fn=<NllLossBackward>)
```

```
 1 model = model_RRRLR()
 2 loss_func = nn.CrossEntropyLoss()
 3 optimizer = optim.Adam(model.parameters(),lr=0.0002)
 4 model.to(device)
 5 loss_array = []
 6 for i in range(num_epoch):
 7     optimizer.zero_grad()
 8     output = model(final_train_tensor).to(device)
```

```
 9     loss = loss_func(output,train_y_tensor)
10     if i%1000==0:
11         print(loss)
12     loss.backward()
13     optimizer.step()
14     loss_array.append(loss)
```

```
tensor(7725928.5000, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(7151.9966, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(31342.8203, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(5864.6211, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(3355.1101, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(1041.5259, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(2898.9993, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(2320.8613, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(813.8175, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(306.7185, device='cuda:0', grad_fn=<NllLossBackward>)
```

```
 1 model = model_RRRR()
 2 loss_func = nn.CrossEntropyLoss()
 3 optimizer = optim.Adam(model.parameters(),lr=0.0005)
 4 model.to(device)
 5 loss_array = []
 6 for i in range(num_epoch):
 7     optimizer.zero_grad()
 8     output = model(final_train_tensor).to(device)
 9     loss = loss_func(output,train_y_tensor)
10     if i%1000==0:
11         print(loss)
12     loss.backward()
13     optimizer.step()
14     loss_array.append(loss)
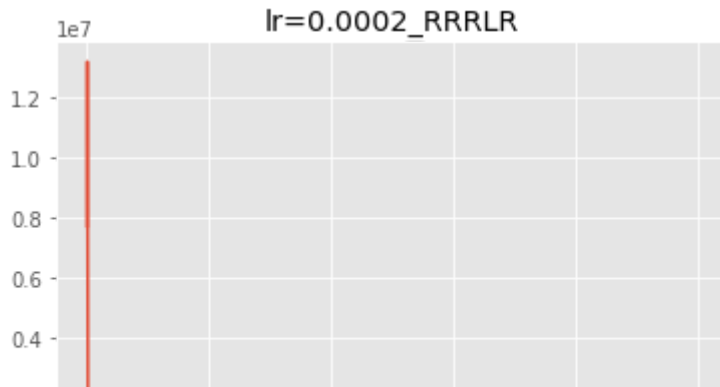```

```
tensor(88838728., device='cuda:0', grad_fn=<NllLossBackward>)
tensor(6955.7974, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(5976.5698, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(7903.1162, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(5816.8589, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(621.4550, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(178.6750, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(140.7881, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(57.2482, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(0.8540, device='cuda:0', grad_fn=<NllLossBackward>)
```

```
 1 # 손실 그래프
 2 plt.plot(loss_array)
 3 plt.title("lr=0.0002_RRRLR")
 4 plt.show()
 5 print(loss_array[-1])
```
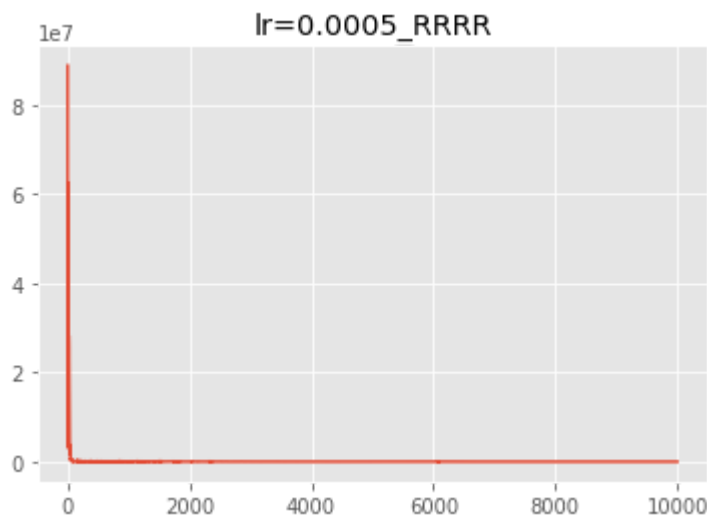
## lr=0.0002_RRRLR

435749  **lr=0.0002_RRRLR.csv**                       2020-06-18  14:12:39    0.842519685

```
1 plt.plot(loss_array)
2 plt.title("lr=0.0005_RRRR")
3 plt.show()
4 print(loss_array[-1])
```

## lr=0.0005_RRRR

tensor(0.1018, device='cuda:0', grad_fn=<NllLossBackward>)

435765  **lr=0.0005_RRRR.csv**                        2020-06-18  15:08:38    0.8346456693

```
1 model = model_TTTT()
2 loss_func = nn.CrossEntropyLoss()
3 optimizer = optim.Adam(model.parameters(),lr=0.005)
4 model.to(device)
5 loss_array = []
6 for i in range(num_epoch):
7     optimizer.zero_grad()
8     output = model(final_train_tensor).to(device)
9     loss = loss_func(output,train_y_tensor)
10    if i%1000==0:
11        print(loss)
12    loss.backward()
13    optimizer.step()
14    loss_array.append(loss)
15 plt.plot(loss_array)
```
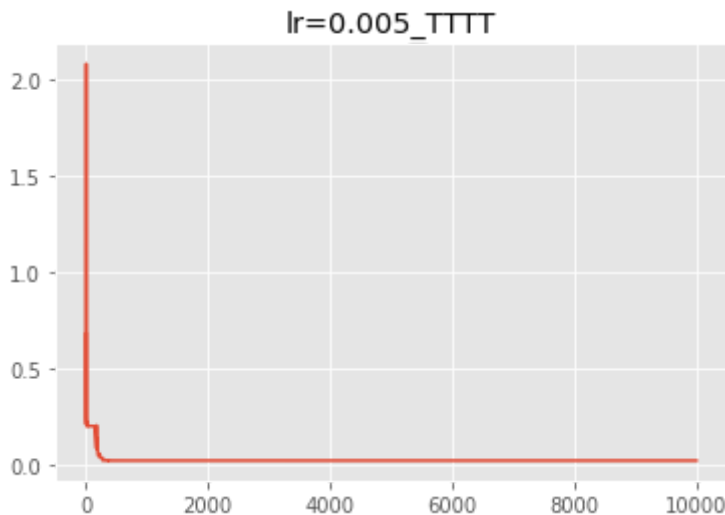
```
16 plt.title("lr=0.005_TTTT")
17 plt.show()
18 print(loss_array[-1])
```

tensor(0.6759, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(0.0185, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(0.0184, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(0.0184, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(0.0184, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(0.0184, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(0.0184, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(0.0184, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(0.0184, device='cuda:0', grad_fn=<NllLossBackward>)
tensor(0.0184, device='cuda:0', grad_fn=<NllLossBackward>)



tensor(0.0184, device='cuda:0', grad_fn=<NllLossBackward>)

---

435766  **lr=0.005_TTTT.csv**                          2020-06-18    0.8503937008
                                                        15:09:01

```
 1 model = model_LRLRLRLR()
 2 loss_func = nn.CrossEntropyLoss()
 3 optimizer = optim.Adam(model.parameters(),lr=0.0002)
 4 model.to(device)
 5 loss_array = []
 6 for i in range(num_epoch):
 7     optimizer.zero_grad()
 8     output = model(final_train_tensor).to(device)
 9     loss = loss_func(output,train_y_tensor)
10     if i%1000==0:
11         print(loss)
12     loss.backward()
13     optimizer.step()
14     loss_array.append(loss)
15
16 plt.plot(loss_array)
17 plt.title("lr=0.0002_LRLRLRLR")
18 plt.show()
19 print(loss_array[-1])
```