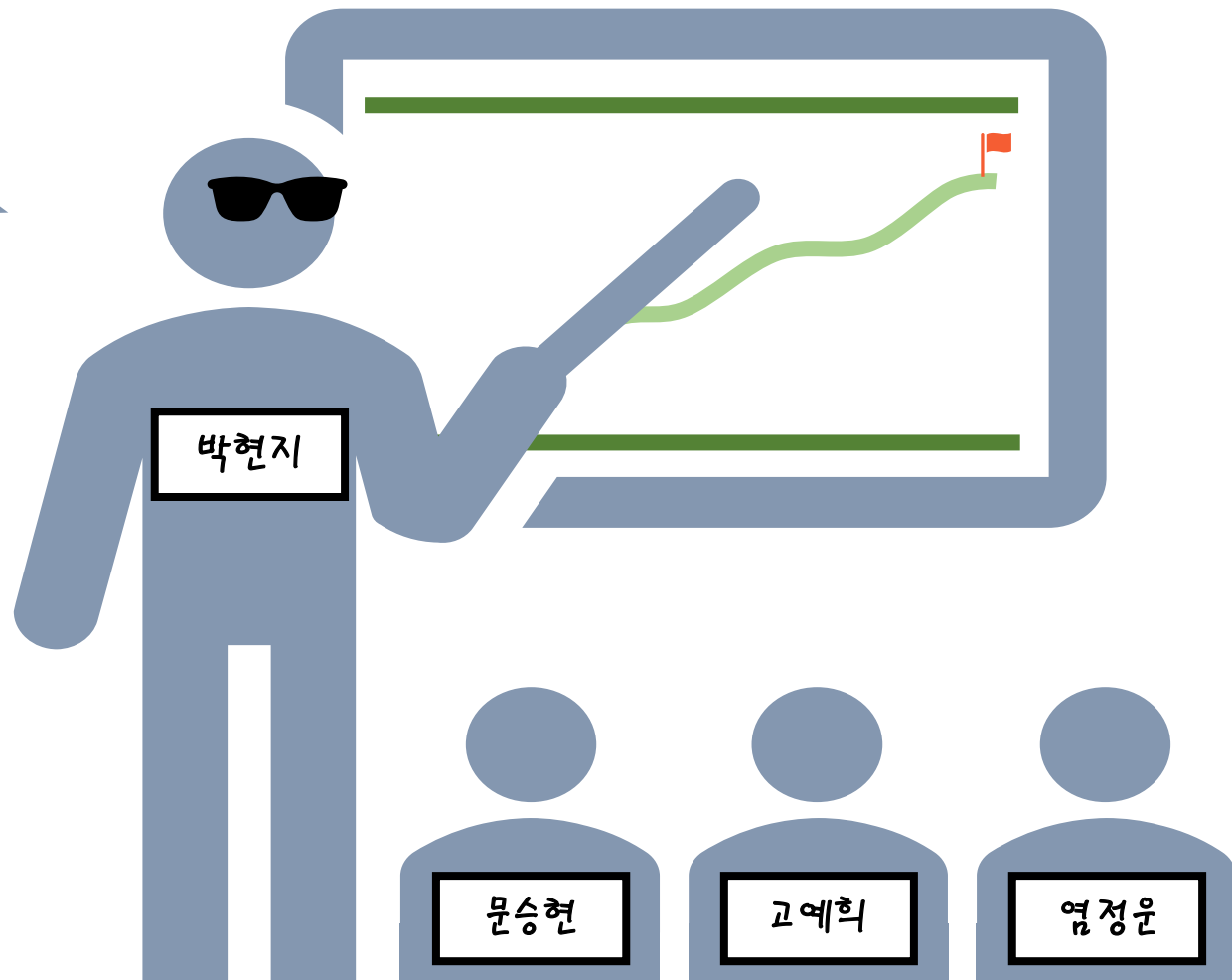
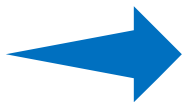


ML4조



목 차



1. 목적



2. 변수분석



3. 결측치처리



4. FE



5. 결과

2주차 - 모델링



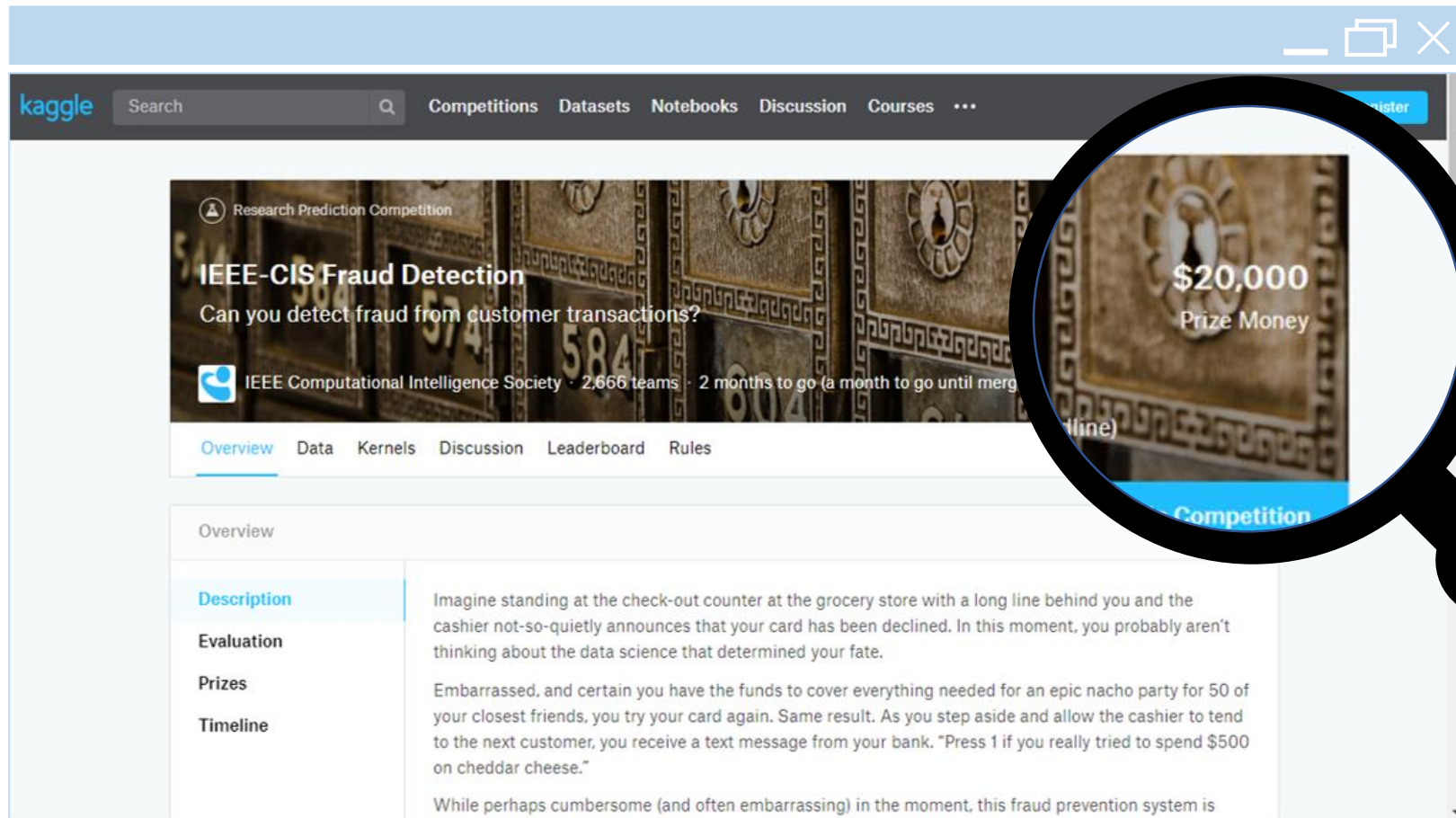
1.1 데이터소개

1.2 변수소개

1.3 최종목표



1.1 데이터소개



The image shows a screenshot of the Kaggle website interface. At the top, the Kaggle logo and navigation links (Search, Competitions, Datasets, Notebooks, Discussion, Courses) are visible. The main content area displays the 'IEEE-CIS Fraud Detection' competition page. A magnifying glass is positioned over the '\$20,000 Prize Money' text, which is highlighted in a blue box. Below the competition title, there is a description of the task: 'Can you detect fraud from customer transactions?'. The page also shows the IEEE Computational Intelligence Society logo, the number of teams (2,666), and the time remaining (2 months to go). The 'Overview' tab is selected, showing a table with columns for Description, Evaluation, Prizes, and Timeline. The Description column contains a paragraph about the competition's context, and the Timeline column contains a paragraph about the fraud prevention system.

Research Prediction Competition

IEEE-CIS Fraud Detection

Can you detect fraud from customer transactions?

IEEE Computational Intelligence Society · 2,666 teams · 2 months to go (a month to go until merge)

Overview Data Kernels Discussion Leaderboard Rules

Overview

Description	Evaluation	Prizes	Timeline
Imagine standing at the check-out counter at the grocery store with a long line behind you and the cashier not-so-quietly announces that your card has been declined. In this moment, you probably aren't thinking about the data science that determined your fate.		Embarrassed, and certain you have the funds to cover everything needed for an epic nacho party for 50 of your closest friends, you try your card again. Same result. As you step aside and allow the cashier to tend to the next customer, you receive a text message from your bank. "Press 1 if you really tried to spend \$500 on cheddar cheese."	While perhaps cumbersome (and often embarrassing) in the moment, this fraud prevention system is

거래데이터
참/거짓 판단

1.2 변수소개

train_transaction.csv



TransactionDT TransactionAMT

ProductCD card1 ~ card6

addr dist P_R_emaildomain

C1-C14 D1-D15 M1-M19 Vxxx

train_identity.csv



DeviceType

DeviceInfo

id1 - id38

test_transaction.csv



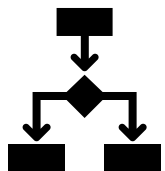
모델 성능 테스트

test_identity.csv

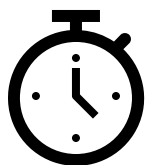


모델 성능 테스트

1.3 최종목표



모델링



빠르고



정확하게



참/거짓



2.1 변수설명

train_transaction.csv

TransactionID TransactionDT

TransactionAmt ProductCD

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain R_emaildomain

C1-C14 D1-D15 M1-M19 Vxxx



TransactionDT

주어진 datetime으로 부터의 timedelta

categorical variable

범주형 변수

continuous variable

연속형 변수

2.1 변수설명

train_transaction.csv



TransactionID TransactionDT

TransactionAmt ProductCD

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain R_emaildomain

C1-C14 D1-D15 M1-M19 Vxxx

categorical variable



범주형 변수

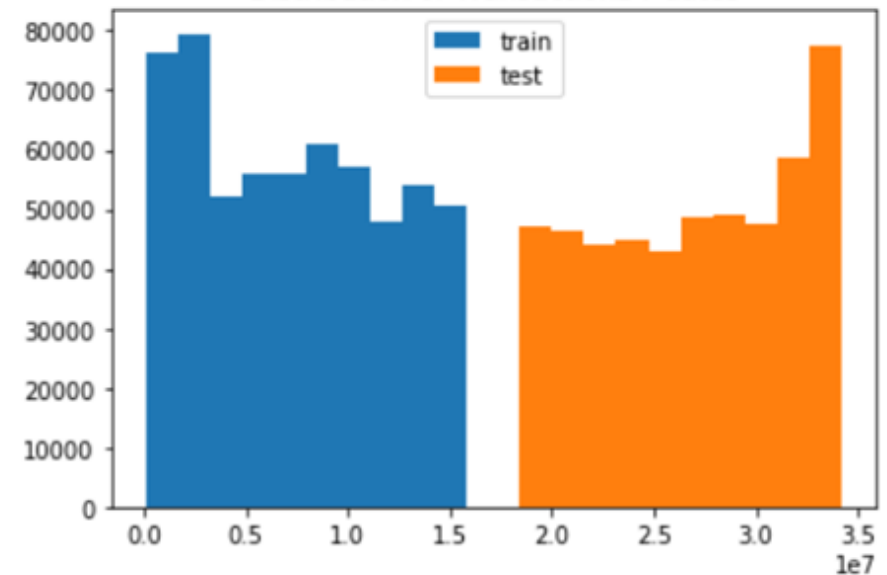
continuous variable



연속형 변수

Text(0.5, 1.0, 'Distribution of TransactionDT dates')

Distribution of TransactionDT dates



Train data와 Test data 사이의 overlapping area가 없다.

2.1 변수설명

train_transaction.csv



TransactionID TransactionDT

TransactionAmt ProductCD

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain R_emaildomain

C1-C14 D1-D15 M1-M19 Vxxx



TransactionAmt



거래 결제 금액(USD)

categorical variable



범주형 변수

continuous variable



연속형 변수

2.1 변수설명

train_transaction.csv

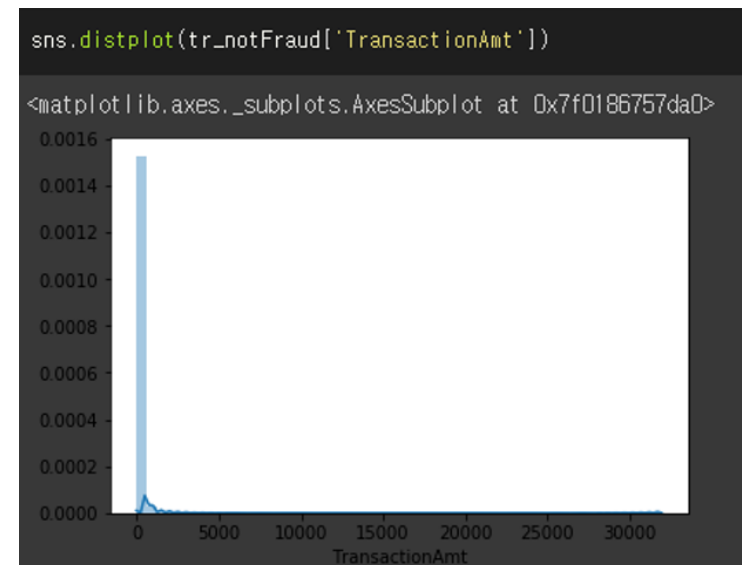
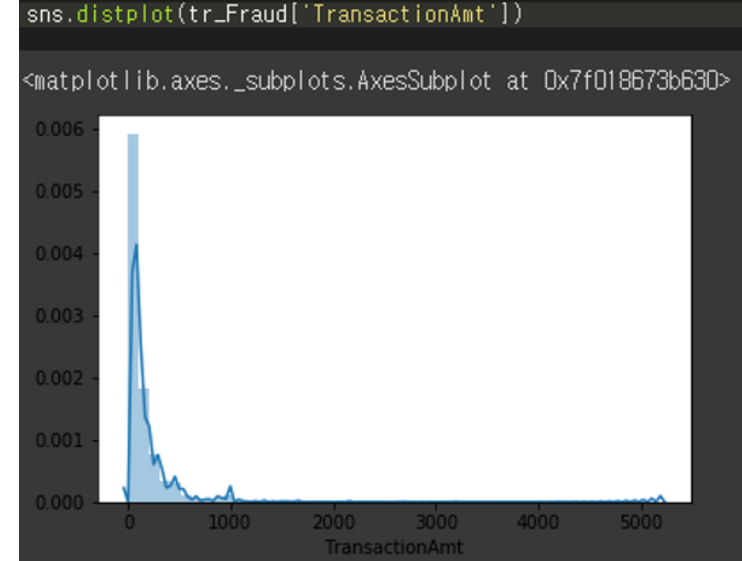
TransactionID		TransactionDT			
TransactionAmt		ProductCD			
card1	card2	card3	card4	card5	card6
Addr1		addr2	dist1		dist2
P_emaildomain		R_emaildomain			
C1-C14	D1-D15	M1-M19		Vxxx	

categorical variable

범주형 변수

continuous variable

연속형 변수



2.1 변수설명

train_transaction.csv



TransactionID TransactionDT

TransactionAmt **ProductCD**

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain R_emaildomain

C1-C14 D1-D15 M1-M19 Vxxx



ProductCD



제품 코드, 각 거래의 제품

categorical variable



범주형 변수

continuous variable



연속형 변수

2.1 변수설명

train_transaction.csv

TransactionID	TransactionDT
TransactionAmt	ProductCD
card1 card2 card3 card4 card5 card6	
Addr1 addr2 dist1 dist2	
P_emaildomain R_emaildomain	
C1-C14 D1-D15 M1-M19 Vxxx	

categorical variable

범주형 변수

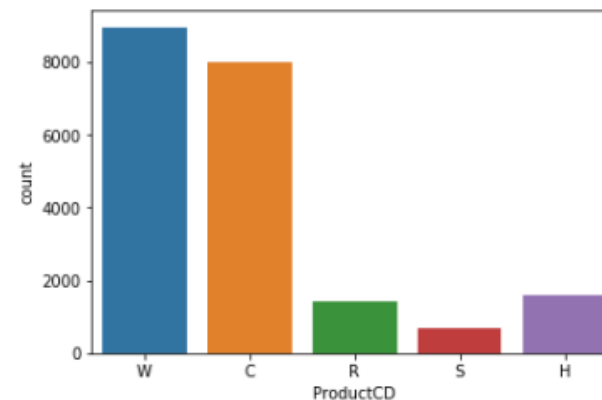
continuous variable

연속형 변수



```
sns.countplot(x = 'ProductCD', data= tr_Fraud, order = ['W', 'C', 'R', 'S', 'H'])  
# plt.title("Fraud Data의 ProductCD")  
# plt.show()
```

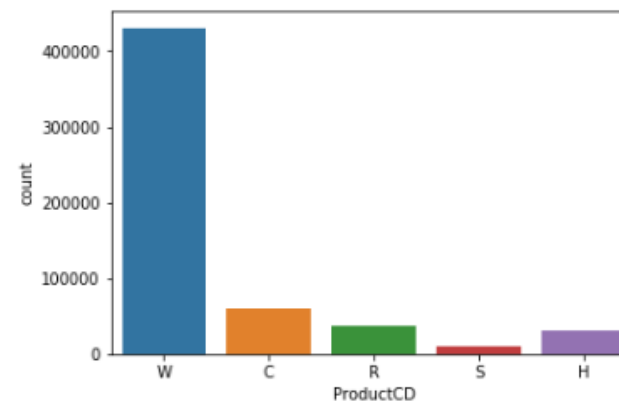
<matplotlib.axes._subplots.AxesSubplot at 0x7fc29a3a9f60>



Fraud

```
[44] sns.countplot(x = 'ProductCD', data= tr_notFraud, order = ['W', 'C', 'R', 'S', 'H'])
```

<matplotlib.axes._subplots.AxesSubplot at 0x7fc29a2e1b00>



Not Fraud

2.1 변수설명

train_transaction.csv

TransactionID TransactionDT

TransactionAmt ProductCD

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain R_emaildomain

C1-C14 D1-D15 M1-M19 Vxxx

categorical variable

범주형 변수

continuous variable

연속형 변수



card1

카드 관련 변수

card2

카드 관련 변수

card3

카드 관련 변수

card4

카드 종류

card5

카드 관련 변수

card6

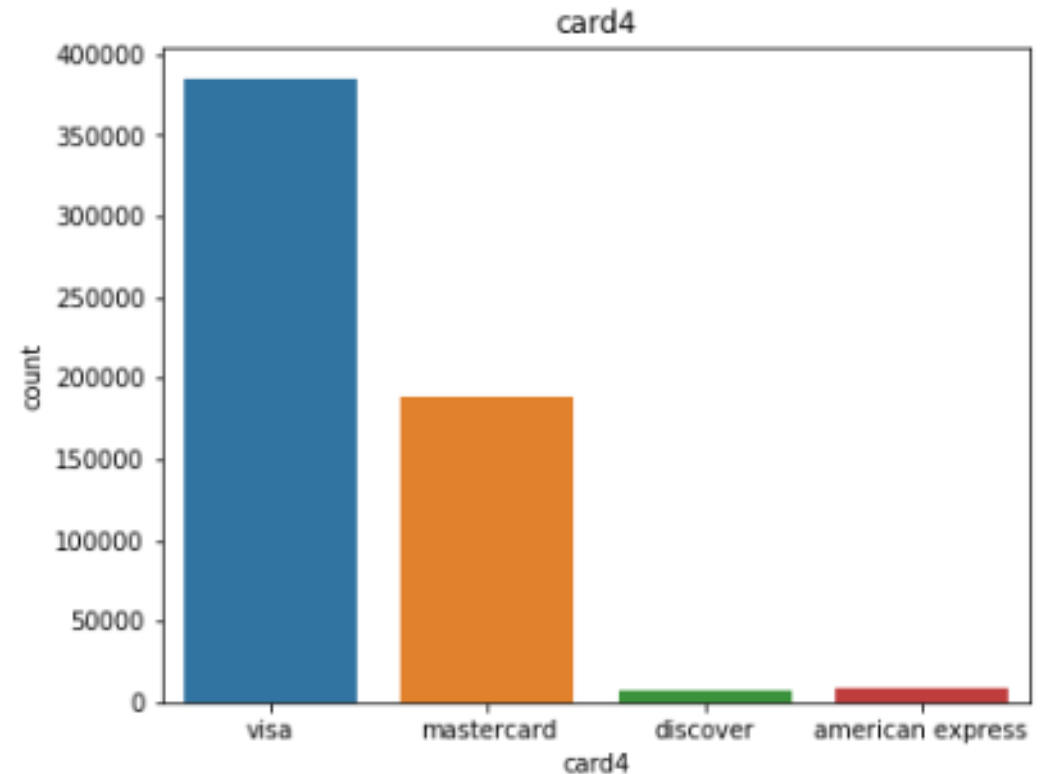
카드 유형

2.1 변수설명

train_transaction.csv



TransactionID TransactionDT
TransactionAmt ProductCD
card1 card2 card3 **card4** card5 card6
Addr1 addr2 dist1 dist2
P_emaildomain R_emaildomain
C1-C14 D1-D15 M1-M19 Vxxx



categorical variable



범주형 변수

continuous variable



연속형 변수

2.1 변수설명

train_transaction.csv

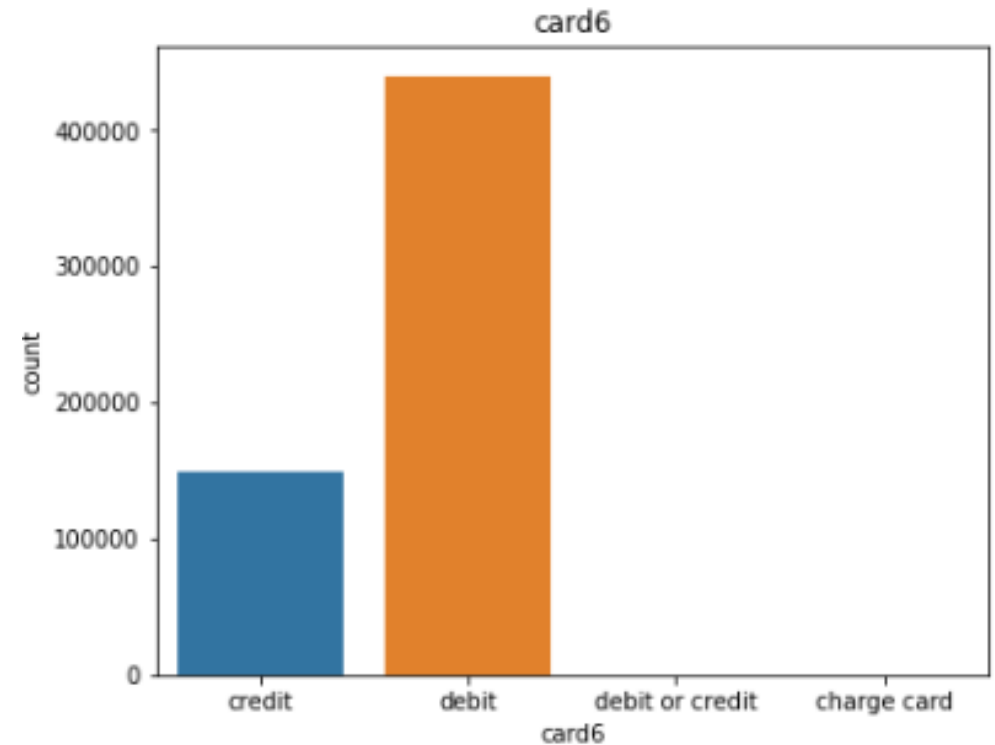
TransactionID TransactionDT
TransactionAmt ProductCD
card1 card2 card3 card4 card5 **card6**
Addr1 addr2 dist1 dist2
P_emaildomain R_emaildomain
C1-C14 D1-D15 M1-M19 Vxxx

categorical variable

범주형 변수

continuous variable

연속형 변수



2.1 변수설명

train_transaction.csv

TransactionID TransactionDT
TransactionAmt ProductCD
card1 card2 card3 card4 card5 card6
Addr1 **addr2** dist1 dist2
P_emaildomain R_emaildomain
C1-C14 D1-D15 M1-M19 Vxxx

categorical variable

범주형 변수

continuous variable

연속형 변수



addr1

도시 주소?

addr2

국가 주소?

2.1 변수설명

train_transaction.csv



TransactionID TransactionDT

TransactionAmt ProductCD

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain R_emaildomain

C1-C14 D1-D15 M1-M19 Vxxx

categorical variable

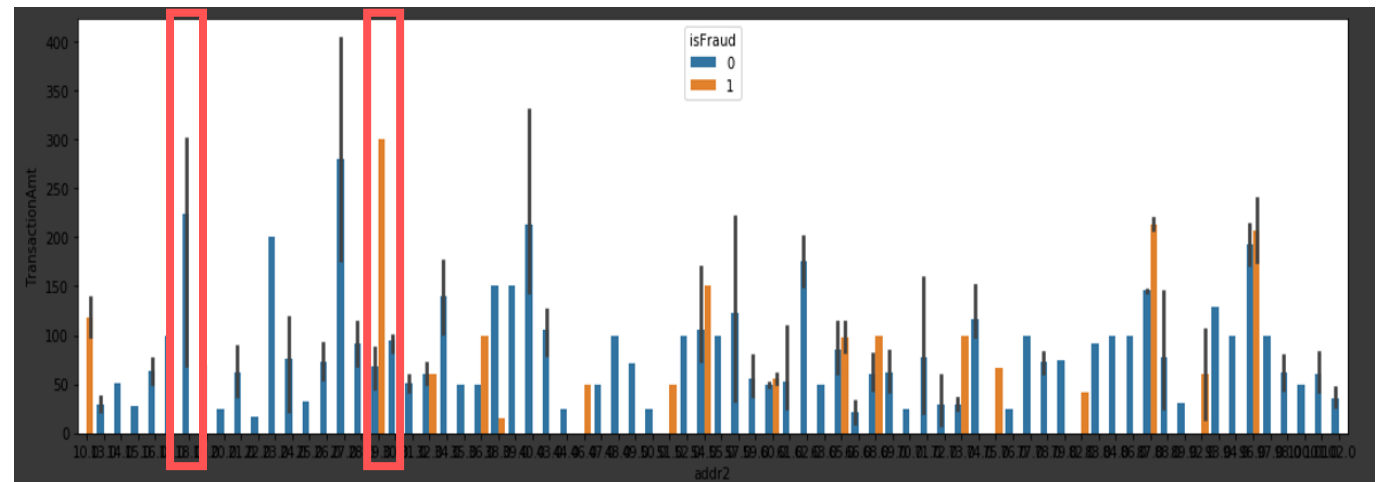


범주형 변수

continuous variable



연속형 변수



특정값에만 데이터가 집중된다.
특정값에서 Fraud / Not Fraud 차이가 크다.

2.1 변수설명

train_transaction.csv

TransactionID TransactionDT
TransactionAmt ProductCD
card1 card2 card3 card4 card5 card6
Addr1 addr2 dist1 dist2
P_emaildomain R_emaildomain
C1-C14 D1-D15 M1-M19 Vxxx

categorical variable

범주형 변수

continuous variable

연속형 변수

dist1

거리 값

dist2

거리 값

2.1 변수설명

train_transaction.csv

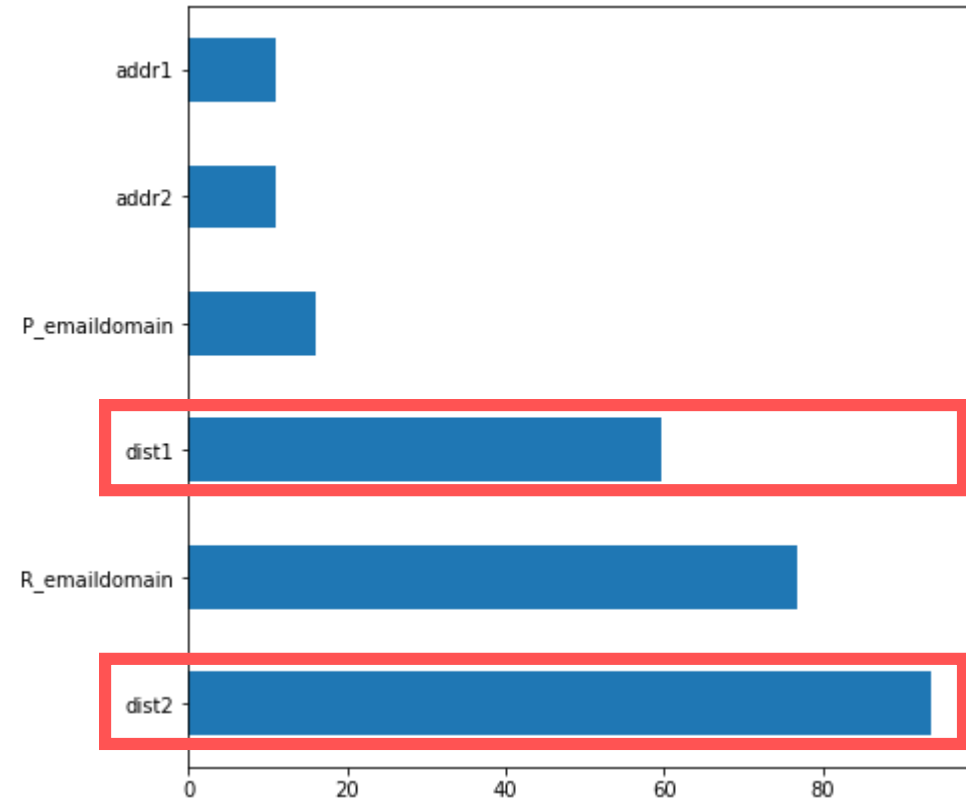
TransactionID	TransactionDT
TransactionAmt	ProductCD
card1 card2 card3 card4 card5 card6	
Addr1 addr2	dist1 dist2
P_emaildomain	R_emaildomain
C1-C14	D1-D15 M1-M19 Vxxx

categorical variable

범주형 변수

continuous variable

연속형 변수



연속형 변수지만 결측치가 너무 많아
Regression이 불가능하기 때문에 DROP 하였다.

2.1 변수설명

train_transaction.csv



TransactionID TransactionDT

TransactionAmt ProductCD

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain R_emaildomain

C1-C14 D1-D15 M1-M19 Vxxx



P_emaildomain



구매자 이메일 도메인

categorical variable



범주형 변수

continuous variable



연속형 변수

2.1 변수설명

train_transaction.csv

TransactionID TransactionDT

TransactionAmt ProductCD

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain R_emaildomain

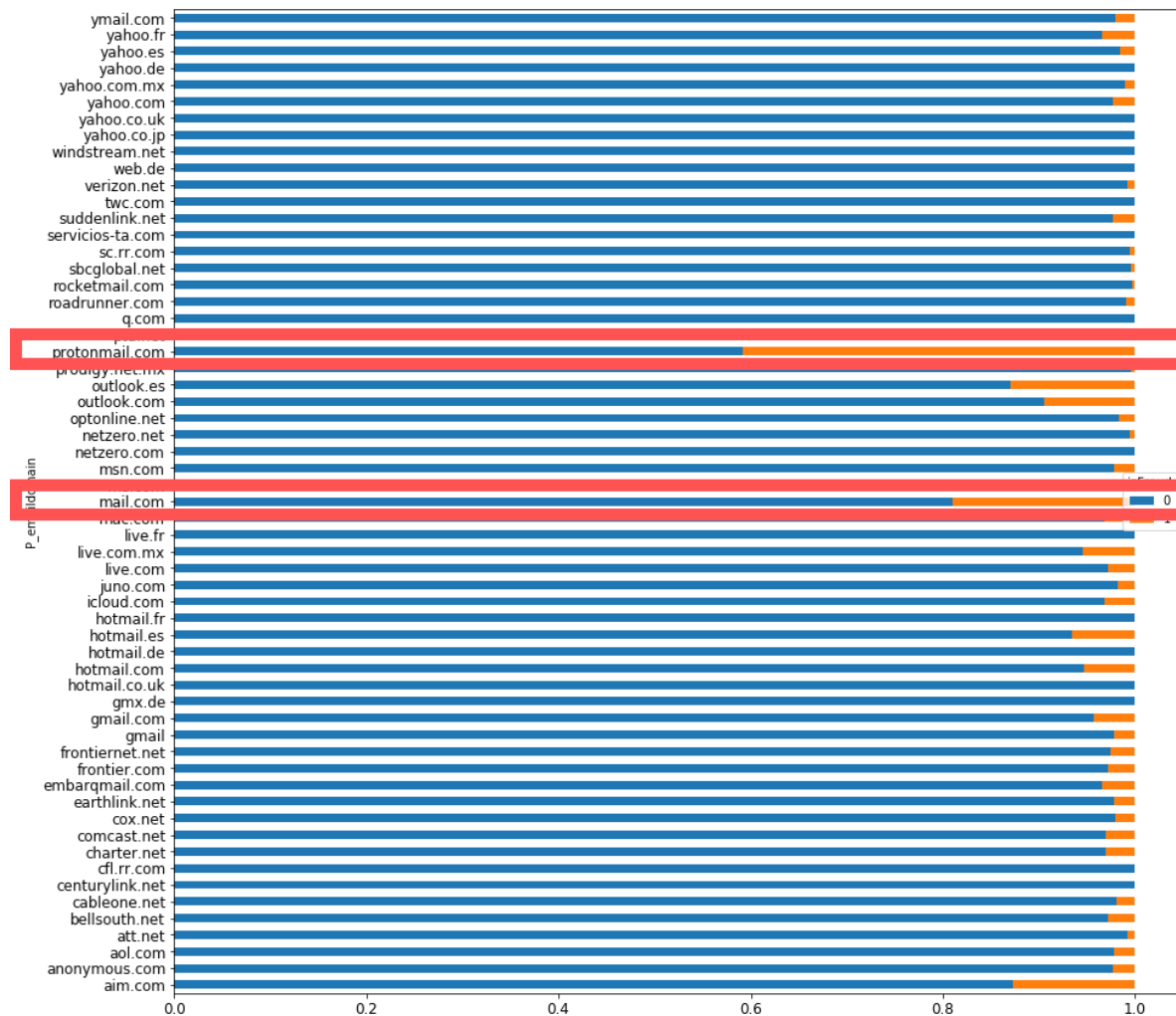
C1-C14 D1-D15 M1-M19 Vxxx

categorical variable

범주형 변수

continuous variable

연속형 변수



이메일에 따라 Fraud / Not Fraud 비율 차이가 크다.

2.1 변수설명

train_transaction.csv



TransactionID TransactionDT

TransactionAmt ProductCD

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain **R_emaildomain**

C1-C14 D1-D15 M1-M19 Vxxx



R_emaildomain



수신자 이메일 도메인

categorical variable



범주형 변수

continuous variable



연속형 변수

2.1 변수설명

train_transaction.csv

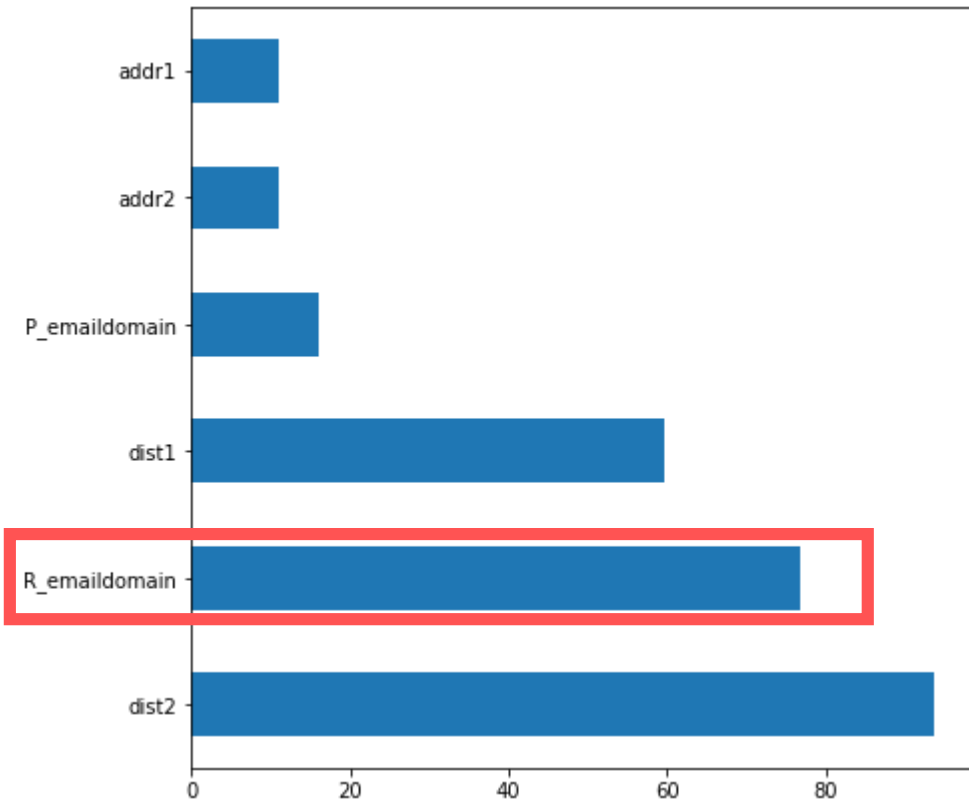
TransactionID	TransactionDT
TransactionAmt	ProductCD
card1 card2 card3 card4 card5 card6	
Addr1 addr2 dist1 dist2	
P_emaildomain	R_emaildomain
C1-C14 D1-D15 M1-M19 Vxxx	

categorical variable

범주형 변수

continuous variable

연속형 변수



NA값이 76.75%로 높기 때문에 DROP 하였다.

2.1 변수설명

train_transaction.csv



TransactionID TransactionDT

TransactionAmt ProductCD

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain R_emaildomain

C1-C14

D1-D15

M1-M19

Vxxx

categorical variable

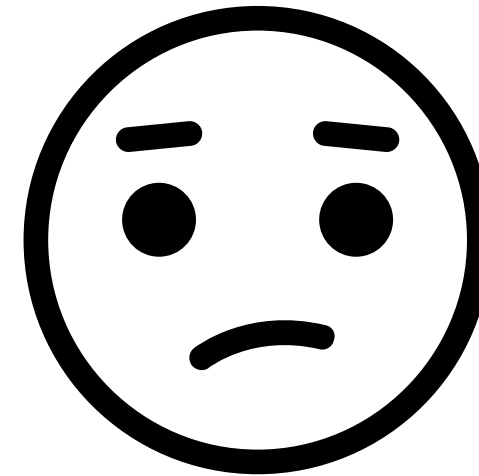


범주형 변수

continuous variable



연속형 변수



의미를 모르는데 변수와 결측치가 너무 많다.

2.1 변수설명

```

-
C1      0.000000
C2      0.000000
C3      0.000000
C4      0.000000
C5      0.000000
C6      0.000000
C7      0.000000
C8      0.000000
C9      0.000000
C10     0.000000
C11     0.000000
C12     0.000000
C13     0.000000
...
V310    0.002032
V311    0.002032
V312    0.002032
V313    0.214888
V314    0.214888
V315    0.214888
V316    0.002032
V317    0.002032
V318    0.002032
V319    0.002032
V320    0.002032
V321    0.002032
V322    86.054967
V323    86.054967
V324    86.054967
V325    86.054967
V326    86.054967
V327    86.054967
V328    86.054967
V329    86.054967
V330    86.054967
V331    86.054967
V332    86.054967
V333    86.054967
V334    86.054967
V335    86.054967
V336    86.054967
V337    86.054967
V338    86.054967
V339    86.054967
Length: 394, dtype: float64

```



```

val = tr_trans[f'V{i}'].isna().sum()/len(tr_trans)+100
if val not in unique_val_V:
    unique_col_V.append(f'V{i}')
    unique_val_V.append(val)

not_needed = []
print(new_tr.shape)
for i in range(1,340):
    if f'V{i}' not in unique_col_V:
        not_needed.append(f'V{i}')

new_tr = new_tr.drop(not_needed, axis = 1)
print(new_tr.shape)

unique_col_C = ['C1']
unique_val_C = []
unique_val_C.append(tr_trans['C1'].isna().sum()/len(tr_trans)+100)

for i in range(2,15):
    val = tr_trans[f'C{i}'].isna().sum()/len(tr_trans)+100
    if val not in unique_val_C:
        unique_col_C.append(f'C{i}')
        unique_val_C.append(val)

not_needed = []
print(new_tr.shape)
for i in range(1,15):
    if f'C{i}' not in unique_col_C:
        not_needed.append(f'C{i}')

new_tr = new_tr.drop(not_needed, axis = 1)
print(new_tr.shape)

unique_col_D = ['D1']
unique_val_D = []
unique_val_D.append(tr_trans['D1'].isna().sum()/len(tr_trans)+100)

for i in range(2,16):
    val = tr_trans[f'D{i}'].isna().sum()/len(tr_trans)+100
    if val not in unique_val_D:
        unique_col_D.append(f'D{i}')
        unique_val_D.append(val)

not_needed = []
print(new_tr.shape)
for i in range(1,16):
    if f'D{i}' not in unique_col_D:
        not_needed.append(f'D{i}')

new_tr = new_tr.drop(not_needed, axis = 1)
print(new_tr.shape)

unique_col_M = ['M1']
unique_val_M = []
unique_val_M.append(tr_trans['M1'].isna().sum()/len(tr_trans)+100)

for i in range(2,10):
    val = tr_trans[f'M{i}'].isna().sum()/len(tr_trans)+100
    if val not in unique_val_M:
        unique_col_M.append(f'M{i}')
        unique_val_M.append(val)

```



```

[ ] card3      0.000000
    card4      0.000000
    card5      0.000000
    card6      0.000000
    addr2      0.000000
    C1         0.000000
    D1         0.214888
    D2         47.549192
    D3         44.514851
    D4         28.604667
    D5         52.467403
    D10        12.873302
    D11        47.293494
    D15        15.090087
    M1         45.907136
    M4         47.658753
    M5         59.349409
    M6         28.678836
    M7         58.635317
    M8         58.633115
    V1         47.293494
    V12        12.881939
    V35        28.612626
    V53        13.055170
    V75        15.098723
    V95        0.053172
    V167       76.355370
    V169       76.323534
    V217       77.913435
    V220       76.053104
    V279       0.002032
    V281       0.214888
    exist      0.000000
    DeviceInfo 0.000000
    dt_month   0.000000
    dt_hour    0.000000
    dtype: float64

```



```

print(added.isna().sum()/len(added))

TransactionID      0.000000
isFraud            0.000000
TransactionDT      0.000000
TransactionAmt     0.000000
ProductCD         0.000000
card1             0.000000
card2             0.000000
card3             0.000000
card4             0.000000
card5             0.000000
card6             0.000000
addr2             0.000000
C1                0.000000
D1                0.214888
D10               12.873302
D15               15.090087
V12               12.881939
V53               13.055170
V75               15.098723
V95               0.053172
V279              0.002032
V281              0.214888
exist             0.000000
DeviceInfo         0.000000
dt_month          0.000000
dt_hour           0.000000
dtype: float64

```

동일 비율
결측치

변수 통합

변수 별
결측치 파악

변수 축약

2.1 변수설명

train_identity.csv	
DeviceType	
DeviceInfo	
id1	id38

categorical variable

범주형 변수

continuous variable

연속형 변수



DeviceType

네트워크 연결 정보(mobile, desktop)

DeviceInfo

네트워크 연결 정보(IOS, Windows)

id1-id38

그 외 네트워크 연결 정보(Browser, IP etc)

2.1 변수설명

train_identity.csv

DeviceType

DeviceInfo

id1 - id38

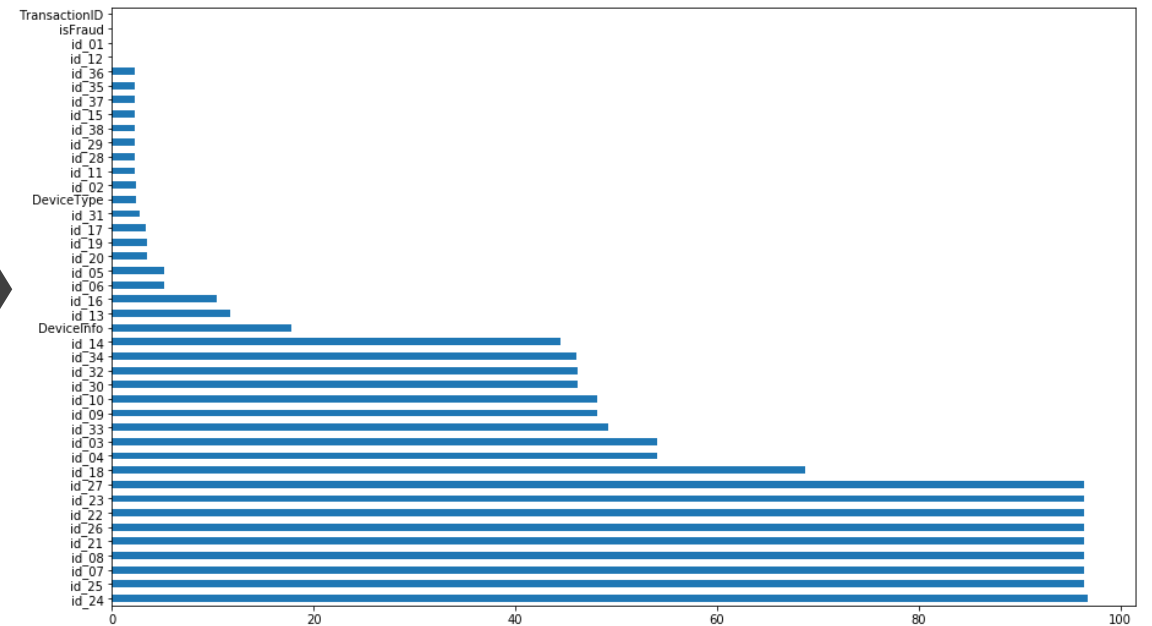
categorical variable

범주형 변수

continuous variable

연속형 변수

Id의 NA값 분포 그래프



2.1 변수설명

train_identity.csv



DeviceType

DeviceInfo

id1 - id38



	TransactionID	isfraud	id_01	id_02	id_03	id_04	id_05	id_06	id_07	id_08	id_09	id_10	id_11	id_13	id_14	id_17	id_18	id_19	id_20	id_21
TransactionID	1	0.0929353	-0.0941453	0.170335	-0.0232686	-0.000267604	-0.0611903	-0.0269559	-0.119807	0.0367722	-0.0329541	0.00924494	0.0310314	0.0850149	0.0431514	0.222829	0.111975	-0.0150304	0.0841049	-0.0290035
isfraud	1	1	-0.120099	0.0493979	0.0414566	-0.0597013	-0.00797841	-0.0271394	-0.0847679	-0.0574083	0.0294309	0.0110426	0.00791442	-0.019536	0.0573237	0.1501	0.0500039	-0.0417214	0.0615965	0.0635439
id_01	-0.0941453	-0.120099	1	-0.142064	0.0195112	0.0396597	0.0096578	0.201637	0.11428	0.095586	0.0291736	0.0507367	0.0144526	0.085596	-0.0653907	-0.180504	-0.0166258	0.000518504	-0.071238	-0.167421
id_02	0.170335	0.0493979	-0.142064	1	-0.0201149	0.000605081	-0.101547	-0.0470409	-0.00075781	0.00504655	-0.0206558	0.01749	0.0537749	-0.0358135	-0.0376841	0.412987	0.142657	-0.0930878	0.108523	-0.00613471
id_03	-0.0232686	0.0414566	0.0195112	-0.0201149	1	0.342178	0.0342799	0.0560694	0.0212344	0.0567367	0.710782	0.112707	-0.00597624	0.0139763	0.00127737	-0.000818863	0.0118563	-0.0095593	0.00334234	0.0718721
id_04	-0.000267604	-0.0597013	0.0396597	0.000605081	0.342178	1	-0.0292644	0.0812691	0.0430196	0.0514552	0.25147	0.337323	-0.00111882	0.010061	-0.0415169	-0.00520691	0.0173001	0.0179257	-0.00660875	0.0111532
id_05	-0.0611903	-0.00797841	0.0096578	-0.101547	0.0342799	-0.0292644	1	-0.291584	0.219281	0.00864465	0.0858817	-0.0709845	-0.0430374	-0.0366375	-0.00198472	-0.203021	-0.137848	-0.0238005	-0.0580697	-0.1019
id_06	-0.0269559	-0.0271394	0.201637	-0.0470409	0.0560694	0.0812691	-0.291584	1	-0.0289097	0.326691	0.0799886	0.222372	0.0208413	0.0748354	-0.0229471	-0.0207075	0.0489046	0.0296467	-0.0710155	0.118837
id_07	-0.119807	-0.0847679	0.11428	0.00075781	0.0212344	0.0430196	0.219281	-0.0289097	1	-0.0940061	0.082746	0.0719875	0.028578	-0.0958749	-0.156789	-0.114993	-0.161544	-0.0509642	-0.0747896	-0.18825
id_08	0.0367722	-0.0574083	0.095586	0.0096578	0.0567367	0.00864465	0.326691	-0.0940061	0.082746	1	0.0863257	-0.0115663	0.0368573	0.00369845	0.0969856	-0.0614565	0.0527358	-0.00183678	0.0855335	0.036373
id_09	-0.0329541	0.0294309	0.0291736	0.00504655	0.710782	0.25147	0.0858817	0.0799886	0.082746	0.112471	1	0.31601	-0.0376715	0.0145817	-0.0066839	-0.0241261	-0.00887793	-0.00615004	0.0436373	0.0363141
id_10	0.00924494	0.0110426	0.0507367	0.0144526	0.112707	0.00597624	0.0700476	0.0223756	0.0863257	0.31601	0.0700476	1	0.0700476	-0.0237586	0.0110145	0.0381123	0.00347172	-0.0304611	0.0181209	-0.141468
id_11	0.0310314	0.00791442	0.0144526	0.085596	-0.0653907	-0.180504	-0.0166258	-0.000518504	-0.071238	-0.167421	-0.0376841	0.412987	1	-0.0632696	0.0170373	0.0626836	0.0273309	0.00404994	0.0363141	-0.020552
id_13	0.0850149	-0.019536	0.085596	-0.0358135	-0.0376841	0.00127737	0.000818863	0.0118563	-0.0095593	0.00334234	0.0718721	-0.071238	-0.167421	1	-0.0810332	-0.0759971	-0.0280719	0.0176573	-0.0294267	0.15894
id_14	0.0431514	0.0573237	-0.0853307	-0.0376841	0.00127737	-0.0415169	-0.00520691	-0.203021	-0.137848	-0.0238005	-0.0580697	-0.1019	-0.0580697	-0.1019	1	-0.086848	0.126425	0.0745749	0.358147	0.0224193
id_17	0.222829	0.1501	-0.180504	0.412987	-0.000818863	0.00520691	-0.203021	-0.137848	-0.0238005	-0.0580697	-0.1019	-0.0580697	-0.1019	-0.086848	0.126425	0.0745749	1	-0.086848	0.126425	0.0745749
id_18	0.111975	-0.0500039	-0.0166258	0.142657	0.0118563	0.0718721	-0.071238	-0.167421	-0.0376841	0.412987	-0.0632696	0.0170373	0.0626836	0.0273309	0.00404994	0.0363141	-0.020552	0.0363141	-0.020552	0.0363141
id_19	-0.0150304	-0.0417214	0.000518504	-0.0930878	-0.0095593	0.0179257	-0.0238005	0.0296467	-0.0509642	0.0527358	-0.00887793	-0.0304611	0.00404994	0.0176573	-0.0745749	-0.214211	0.0197265	1	-0.0883086	0.253749
id_20	0.0841049	0.0615965	-0.071238	0.108523	0.00334234	-0.00660875	-0.0580697	-0.1019	-0.0580697	-0.1019	-0.0580697	-0.1019	-0.0580697	-0.1019	-0.086848	0.126425	0.0745749	0.358147	0.0224193	0.0363141
id_21	-0.0290035	0.0635439	-0.167421	-0.00613471	0.0718721	0.0111532	-0.1019	0.118837	-0.18825	0.0855335	0.0436373	-0.141468	-0.100552	0.15894	0.0224193	0.0432042	0.202953	0.253749	0.0010873	1
id_22	0.0526215	0.118409	0.00824133	0.190991	0.0686015	0.000459111	-0.126581	0.131131	-0.277448	0.143301	0.0322286	0.0504429	-0.0207672	0.0965755	0.100784	0.591957	0.130109	0.0638709	0.0307648	0.0705909
id_24	-0.0383391	-0.00190505	-0.132626	0.0397599	-0.0227825	-0.0731104	-0.00658385	0.0752008	-0.070752	-0.0118495	-0.0123118	-0.12249	-0.0956444	0.0682135	-0.111833	-0.0148953	0.014321	0.245221	-0.189828	0.220933
id_25	0.0206723	0.0340447	-0.0369538	0.0237625	-0.0182474	-0.0849645	-0.0574632	-0.0523055	0.0376495	-0.00832803	0.00312764	0.0021228	0.0539387	0.483726	-0.0504967	0.173535	0.108161	-0.00805921	-0.0284514	-0.147494
id_26	0.0136696	0.099537	-0.0698744	0.0690952	0.025608	-0.035366	-0.0336025	0.067236	-0.131638	0.0372123	-0.0540884	-0.0258323	-0.0661955	0.00736990	0.116977	0.076164	-0.0487985	0.116386	0.0484128	0.0537209
id_32	-0.0664368	0.0697017	-0.0438616	0.130711	-0.0115261	0.0164673	0.0996038	-0.0863028	0.227533	-0.0023808	-0.0295331	0.0568071	0.152434	-0.0208539	-0.0195325	0.0175765	0.069	-0.0842104	0.0035675	-0.431019
id count	0.0896217	0.0120222	-0.070876	0.373101	0.0138708	0.0248379	-0.2124	0.0906691	-0.207424	0.156631	0.00409214	0.0717117	0.0199941	-0.299455	-0.0938074	0.662938	0.282309	-0.124806	0.160293	0.2588

categorical variable

범주형 변수

continuous variable

연속형 변수

id3과 id9의 상관관계가 높다.



1. 목적



2. 변수분석



3. 결측치처리





3.1 방법소개

3.2 결측치비율

3.3 결측치처리

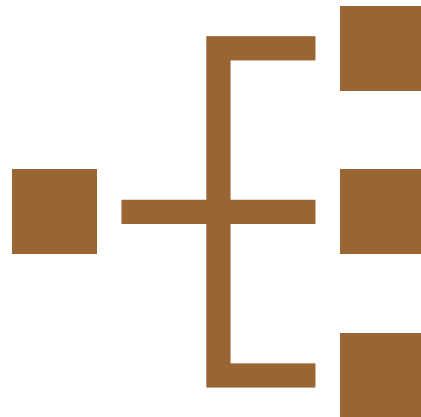


3.1 방법소개

결측값처리

고려사항

- 변수간 관계를 파악
- 왜곡을 적게 만들어 모델의 정확성 높임
- 결측값이 무작위인가? 관계가 있는가?



A. 삭제

결측치가 생긴부분을 삭제

장점 : 빠르고 간편

단점 : 모델의 유효성이 낮아짐 무작위발생이 아닐 시 왜곡된 모델 생성

B. 대체 (평균, 최빈값, 중간값)

일괄 대체 / 범주형 변수로 유사유형 평균값으로 대체

장점 : 빠르고 간편

단점 : 모델의 유효성이 낮아짐, 유사유형 선택 시 왜곡된 모델 생성

C. 예측값 삽입

결측값이 없는 관측치로 예측 모델 생성 (Regression, Logistic regression 등)

장점 : 자의적인 판단이 적음

단점 : 다양한 변수에서 결측치가 발생하거나 결측치가 많은 경우 사용불가

3.2 결측치비율

	전체	결측치 수	결측치 비율(%)
TransactionID	590541	0	0
isFraud	590541	0	0
TransactionDT	590541	0	0
TransactionAmt	590541	0	0
ProductCD	590541	0	0
card1	590541	0	0
card2	590541	8933	1.51
card3	590541	1565	0.26
card4	590541	1577	0.27
card5	590541	4259	0.72
card6	590541	1571	0.27
addr1	590541	65706	11.1
addr2	590541	65706	11.1
dist1	590541	352271	59.65
dist2	590541	552913	93.62
P_emaildomain	590541	94456	15.99
R_emaildomain	590541	453249	76.75

train_transaction.csv



TransactionDT TransactionAMT

ProductCD card1 ~ card6

addr dist P_R_emaildomain

C1-C14 D1-D15 M1-M19 Vxxx

3.2 결측치비율

dist1



결측치 59.652%

dist2

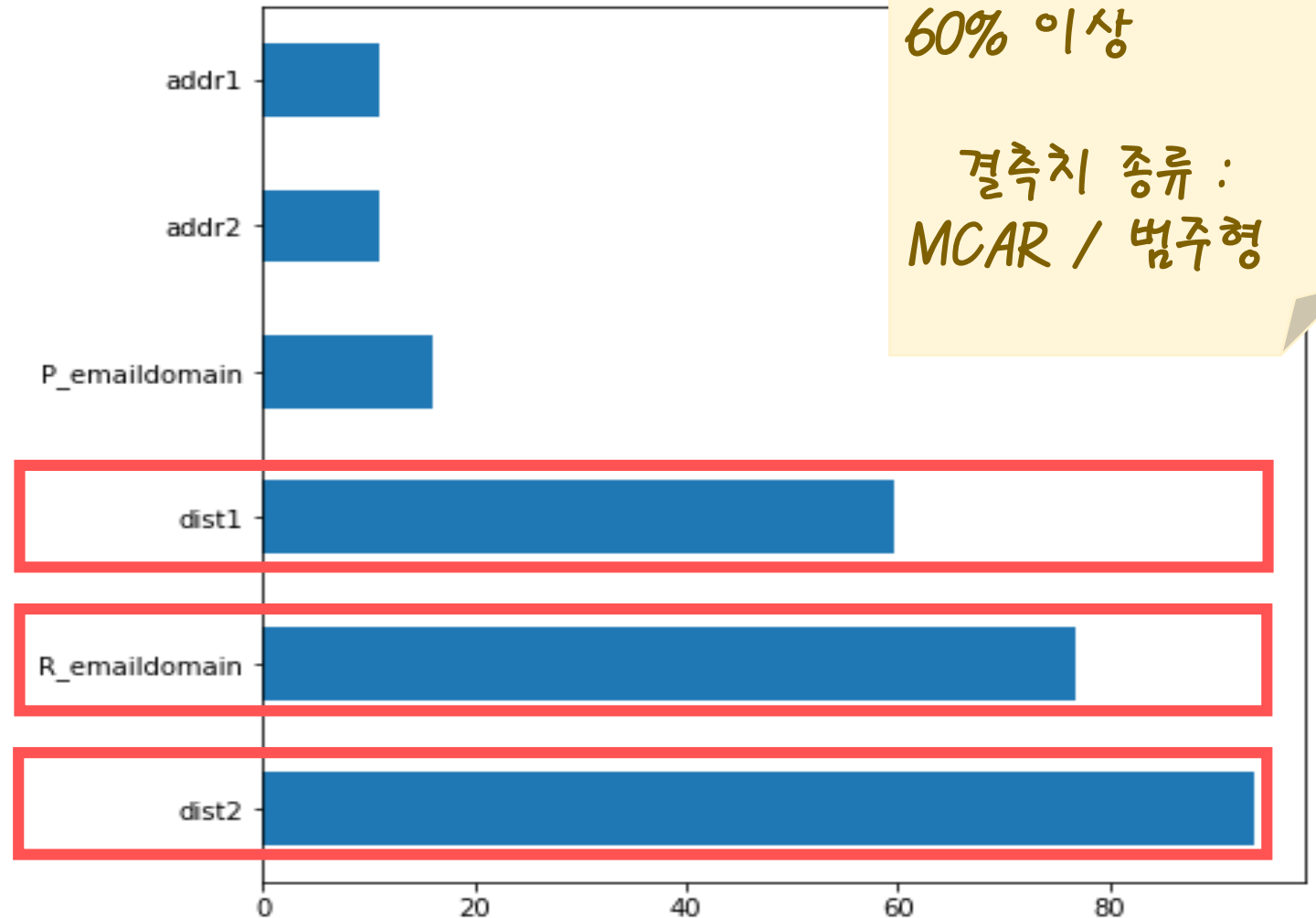


결측치 93.628%

R_emaildomain



결측치 76.751%



3.2 결측치비율

addr1

도시

?



addr2

국가

?



addr1 & addr2

결측치가 완벽하게
일치하므로 하나로 통합

3.2 결측치비율

1	Transac	isFraud	Fraud	Transac	Transac	Product	card1	card2	card3	card4	card5	card6	addr1	addr2	dist1	dist2	P_email	R_email	
65	2987063	0	notFraud	87604	80	W	17967	285	150	visa		226	debit	184	87			yahoo.com	
66	2987064	0	notFraud	87611	250	W	3278	453	150	visa		226	debit	122	87			gmail.com	
67	2987065	0	notFraud	87650	114.95	W	17188	321	150	visa		226	debit	299	87	1			
68	2987066	0	notFraud	87660	300	H	15333	562	150	visa		226	credit	315	87			anonymous.co	
69	2987067	0	notFraud	87664	445	W	1240	302	150	visa		195	credit	485	87	24			
70	2987068	0	notFraud	87667	3.081	C	14076	545	185	visa		147	credit					hotmail.cc hoti	
71	2987069	0	notFraud	87725	20	S	12866	303	150	visa		226	debit	330	87		84		hoti
72	2987070	0	notFraud	87735	100	H	3682	264	150	visa		162	credit	325	87			anonymous.co	
73	2987071	0	notFraud	87736	59	W	1662	111	150	visa		195	debit	472	87			gmail.com	
74	2987072	0	notFraud	87752	6.767	C	13832	375	185	mastercard		224	debit					outlook.cc outl	
75	2987073	0	notFraud	87759	554	W	1955	383	150	visa		226	debit	315	87	5		yahoo.com	
76	2987074	0	notFraud	87775	27.793	C	15885	545	185	visa		138	debit				100	gmail.com gma	
77	2987075	0	notFraud	87779	68.5	W	4806	490	150	visa		226	debit	315	87			gmail.com	
78	2987076	0	notFraud	87787	36.95	W	18132	567	150	mastercard		117	debit	441	87	1			
79	2987077	0	notFraud	87793	280	W	3278	453	150	visa		226	debit	122	87			gmail.com	
80	2987078	0	notFraud	87825	300	W	14858	558	150	visa		226	debit	220	87	2239			
81	2987079	0	notFraud	87839	28.699	C	4504	500	185	mastercard		219	credit					hotmail.cc hoti	
82	2987080	0	notFraud	87856	60	W	7207	111	150	visa		226	debit	204	87			comcast.net	
83	2987081	0	notFraud	87868	104.95	W	17188	321	150	visa		226	debit	299	87	1		gmail.com	
84	2987082	0	notFraud	87899	280	W	15066	170	150	mastercard		102	credit	325	87			optonline.net	
85	2987083	0	notFraud	87924	411.95	W	14695	396	150	mastercard		224	credit	315	87	22		gmail.com	
86	2987084	0	notFraud	87928	125.674	C	5583	103	185	visa		226	credit				744	anonymou ano	
87	2987085	0	notFraud	87935	42.596	C	15885	545	185	visa		138	debit					anonymou ano	
88	2987086	0	notFraud	87950	44.5	W	11815	206	150	mastercard		166	debit	476	87	6		gmail.com	
89	2987087	0	notFraud	88004	88.95	W	1549	143	150	visa		226	debit	205	87	9		yahoo.com	
90	2987088	0	notFraud	88021	140	W	18095	243	150	visa		226	debit	387	87			gmail.com	
91	2987089	0	notFraud	88042	318.95	W	17188	321	150	visa		226	debit	299	87	5		gmail.com	
92	2987090	0	notFraud	88046	77.821	C	15885	545	185	visa		138	debit					gmail.com gma	
93	2987091	0	notFraud	88053	107	W	14165	111	150	mastercard		224	debit	181	87	10			
94	2987092	0	notFraud	88054	117	W	6481	111	150	visa		226	debit	337	87	327		cox.net	
95	2987093	0	notFraud	88070	50	H	5220	360	150	visa		226	credit	231	87			charter.net cha	
96	2987094	0	notFraud	88107	527	W	5409	170	150	visa		226	credit	315	87	5		yahoo.com	
97	2987095	0	notFraud	88120	59	W	2538	476	150	visa		166	debit	330	87			anonymous.co	

3.2 결측치비율

	전체	결측치 수	결측치 비율(%)
TransactionID	590541	0	0
isFraud	590541	0	0
TransactionDT	590541	0	0
TransactionAmt	590541	0	0
ProductCD	590541	0	0
card1	590541	0	0
card2	590541	8933	1.51
card3	590541	1565	0.26
card4	590541	1577	0.27
card5	590541	4259	0.72
card6	590541	1571	0.27
addr1	590541	65706	11.1
addr2	590541	65706	11.1
dist1	590541	352271	59.65
dist2	590541	552913	93.62
P_emaildomain	590541	94456	15.99
R_emaildomain	590541	453249	76.75



	전체	결측치 수	결측치 비율(%)
TransactionID	590541	0	0
isFraud	590541	0	0
TransactionDT	590541	0	0
TransactionAmt	590541	0	0
ProductCD	590541	0	0
card1	590541	0	0
card2	590541	8933	1.51
card3	590541	1565	1.51
card4	590541	1577	0.26
card5	590541	4259	0.27
card6	590541	1571	0.72
addr2	590541	65706	11.1
P_emaildomain	590541	94456	15.99

3.3 결측치처리

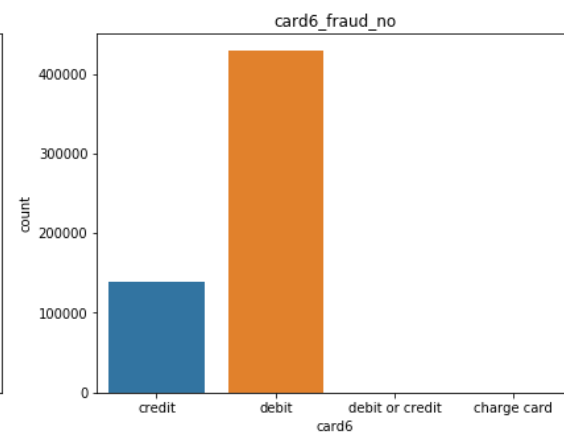
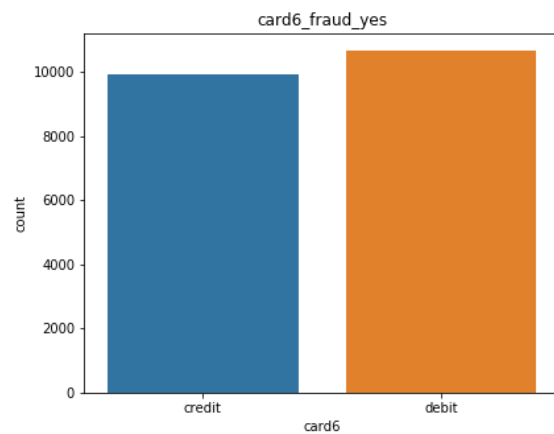
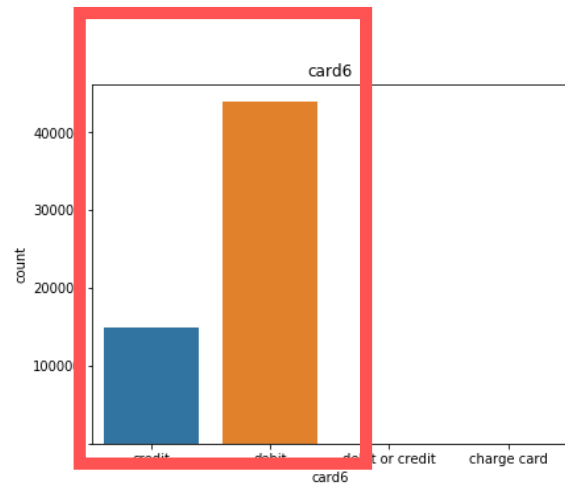
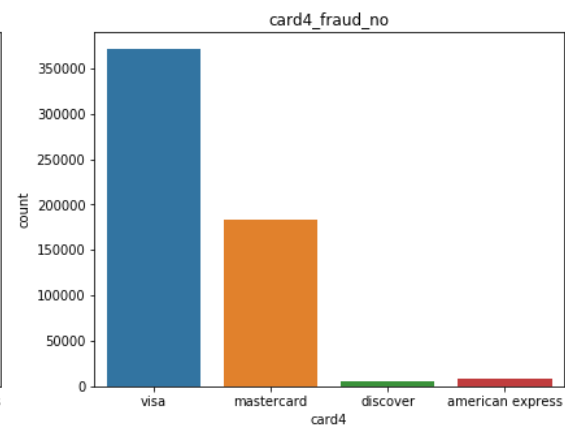
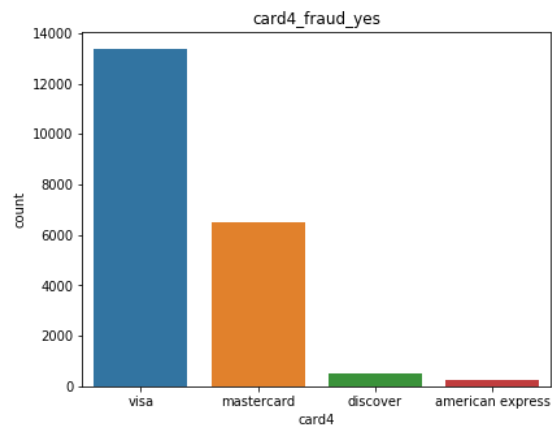
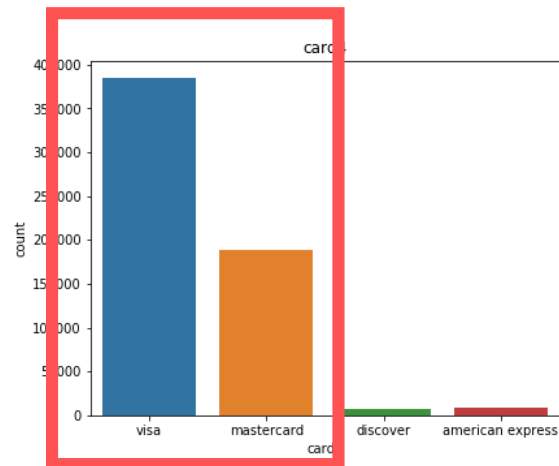
Card4

Card6

최빈값으로 대체

결측치 비율이 낮으며, 범주가 많지 않아 최빈값 대체가능

3.3 결측치처리



3.3 결측치처리

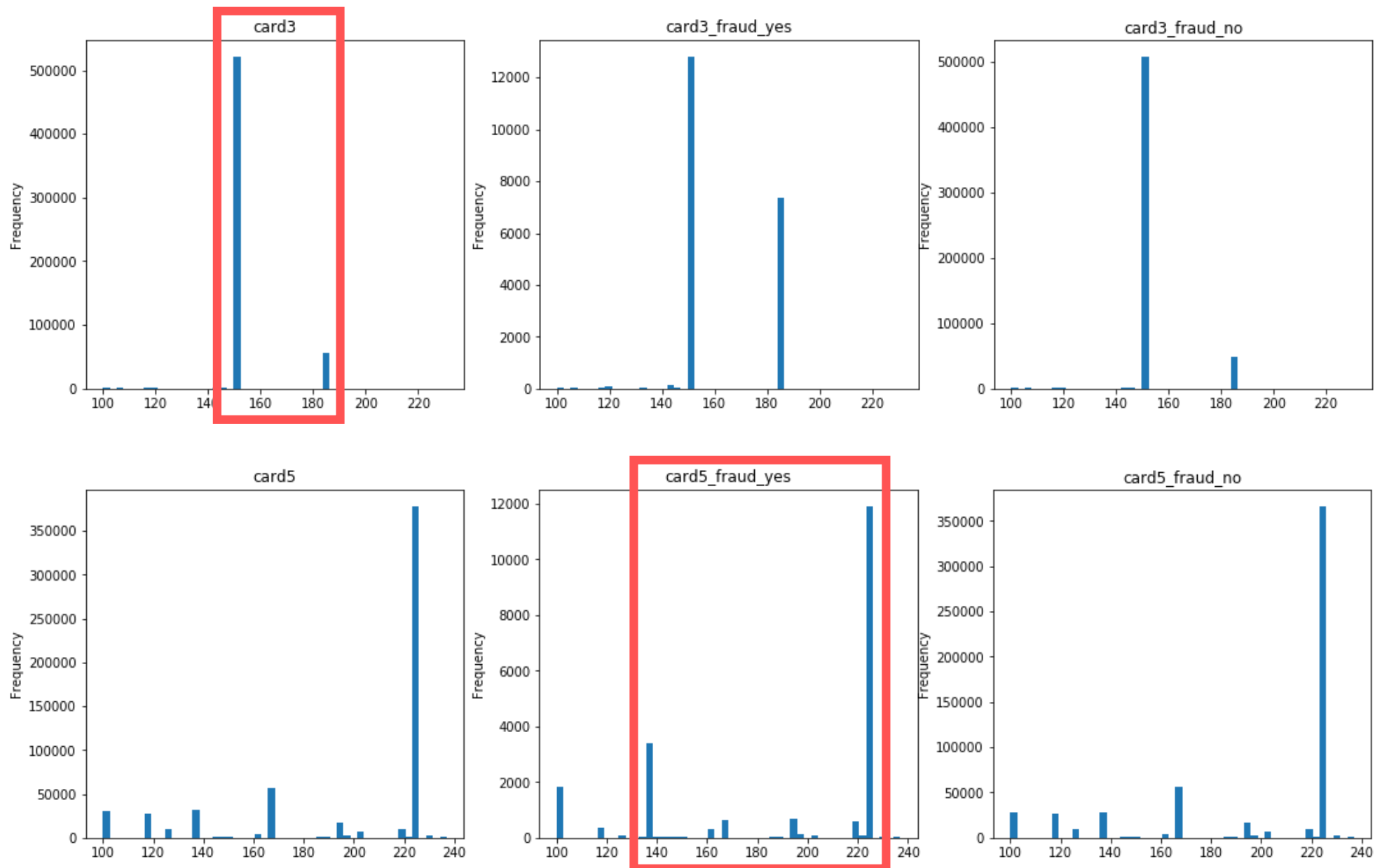
Card3

Card5

최빈값으로 대체

그래프에서 적은 수의 변수만 유의미한 수치로 나옴

3.3 결측치처리



3.3 결측치처리

DeviceInfo

Unknown 처리

현실적으로 기기의 기종을 찾기 불가 (MCAR)

특히, 기기의 기종, OS이름, 브라우저 엔진 등

여러 다른 변수가 섞여있는 Column

Outlier 처리?

이상치(Outlier)의 경우
Fraud Detection에서는
주요 요인이 될 수 있음

이상치 처리에 영향을
크게 받지 않는

Decision Tree 기반 모델 이용 예정



1. 목적



2. 변수분석

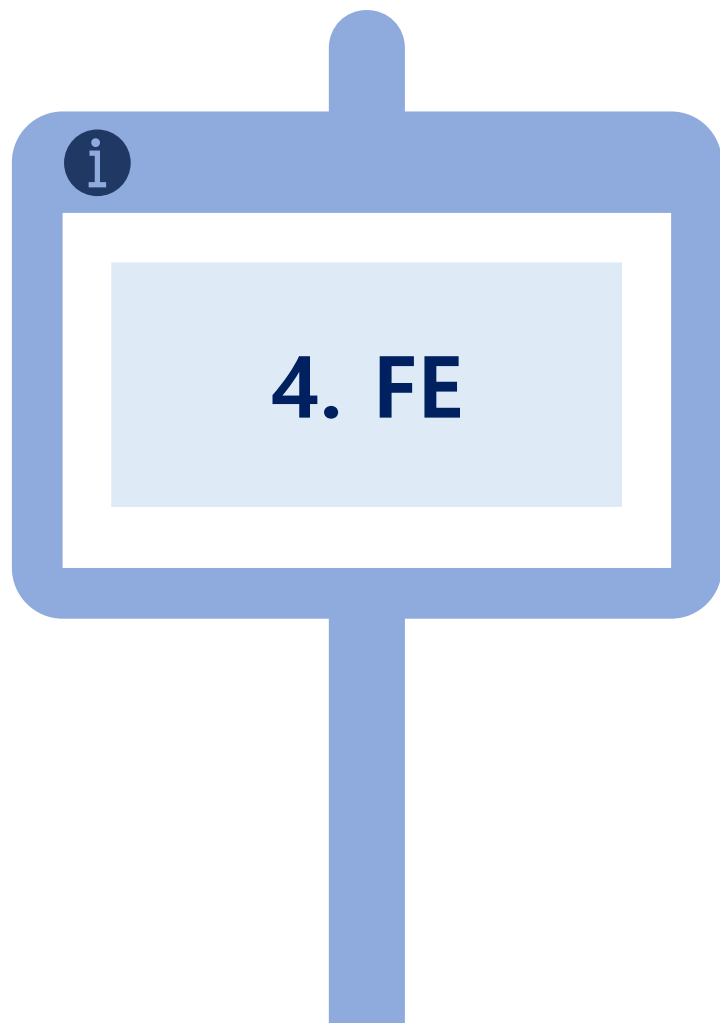


3. 결측치처리



4. FE





4.1 FE소개

4.2 FE결과



Feature Engineering

기존 변수 사용 -> 데이터 추가

관측치 / 변수 추가 없이 기존 데이터 강화

- Scaling -> log / root
- Binning -> 연속형을 범주형으로
- Transform -> 새로운 변수 추가 및 결합
- Dummy -> One-Hot Encoding

From Ybigta 세션자료

The machine learning models affected by the magnitude of the feature are:

- Linear and Logistic Regression
- Neural Networks
- Support Vector Machines
- KNN
- K-means clustering
- Linear Discriminant Analysis (LDA)
- Principal Component Analysis (PCA)

*Machine learning models **insensitive to feature magnitude** are the ones based on Trees:*

- Classification and Regression Trees
- Random Forests
- Gradient Boosted Trees

Feature Engineering

기존

관측치 / 변수 추가 없이 기존 데이터 강화

- Scaling -> 정규화
- Binning -> 연속형을 범주형으로
- Transform -> 새로운 변수 추가 및 결합
- Dummy -> One-Hot Encoding

From Ybigta 세션자료

The machine learning models affected by the magnitude of the feature are:

• Logistic Regression
• Neural Networks
• Support Vector Machines

• K-means clustering

• Linear Discriminant Analysis (LDA)

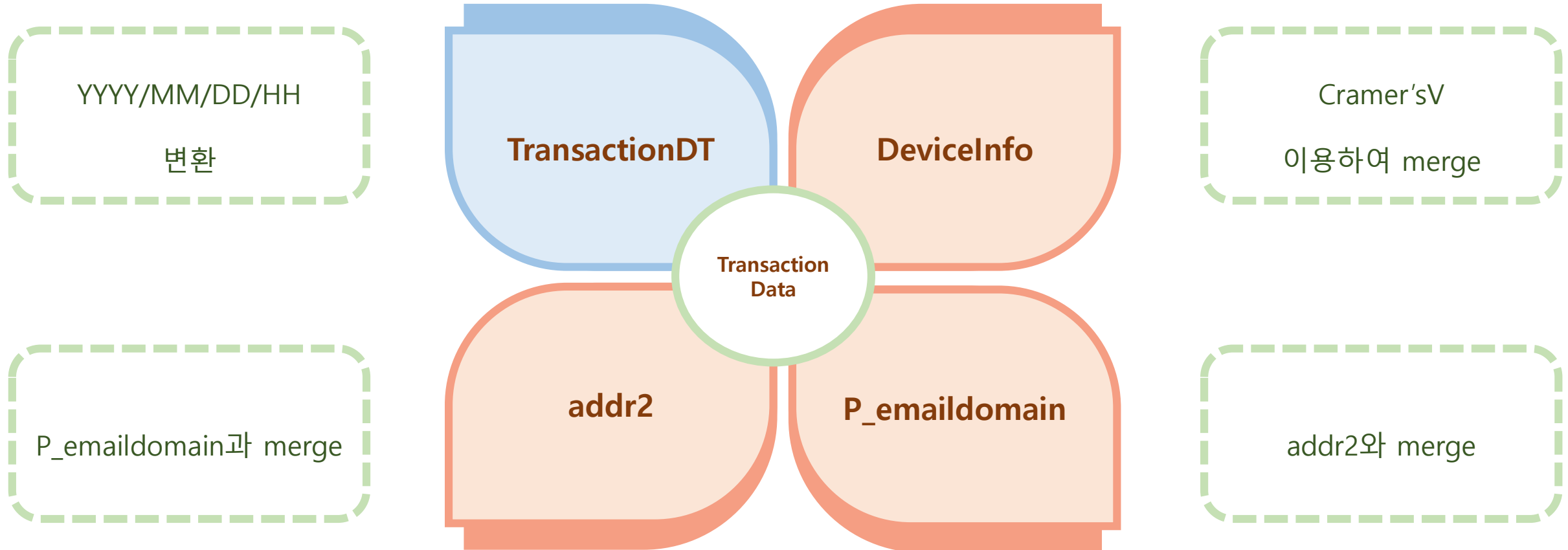
• Principal Component Analysis (PCA)

Machine learning models insensitive to feature magnitude are the ones based on Trees:

- Classification and Regression Trees
- Random Forests
- Gradient Boosted Trees

Gradient Boost계열 모델 이용 예정
Magnitude feature 재조정 필요 없음

4.2 FE결과



4.2 FE결과

TransactionDT

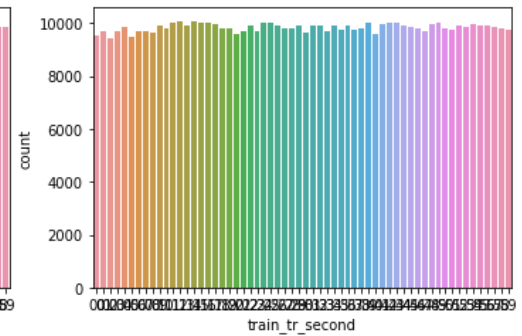
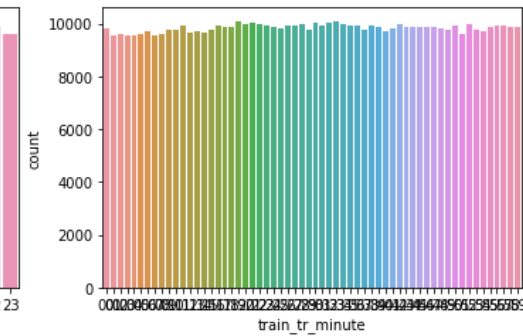
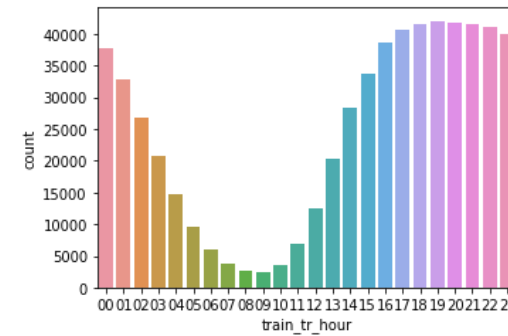
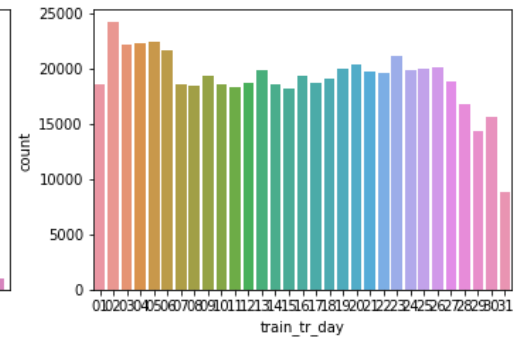
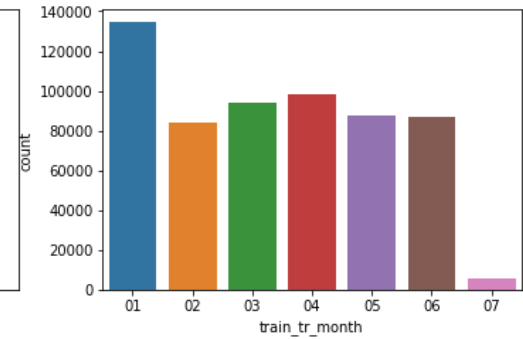
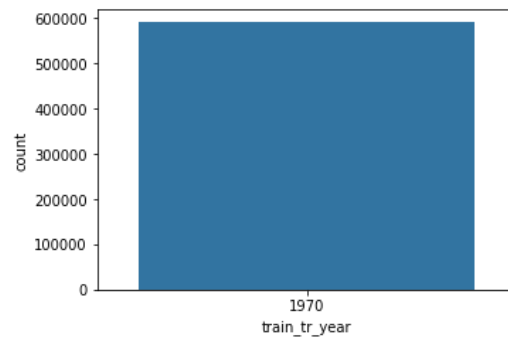
```
train_tr.time["train_tr_year"] = train_tr.time["TransactionDT_cleaned"].astype(str).str[4]
train_tr.time["train_tr_month"] = train_tr.time["TransactionDT_cleaned"].astype(str).str[5:7]
train_tr.time["train_tr_day"] = train_tr.time["TransactionDT_cleaned"].astype(str).str[8:10]
train_tr.time["train_tr_hour"] = train_tr.time["TransactionDT_cleaned"].astype(str).str[11:13]
train_tr.time["train_tr_minute"] = train_tr.time["TransactionDT_cleaned"].astype(str).str[14:16]
train_tr.time["train_tr_second"] = train_tr.time["TransactionDT_cleaned"].astype(str).str[17:]

train_tr.time["train_tr_dayofweek"] = train_tr.time["TransactionDT_cleaned"].dt.dayofweek
print(train_tr.time.shape)
train_tr.time[["train_tr_year", "train_tr_month", "train_tr_day", "train_tr_hour", "train_tr_minute", "train_tr_second", "train_tr_dayofweek"]].head()
(3540, 402)
```

train_tr_year	train_tr_month	train_tr_day	train_tr_hour	train_tr_minute	train_tr_second	train_tr_dayofweek
1970	01	02	00	00	00	4
1970	01	02	00	00	01	4
1970	01	02	00	01	09	4
1970	01	02	00	01	39	4
1970	01	02	00	01	46	4

```
train_tr.time.loc[train_tr.time["train_tr_dayofweek"] == 0, "train_tr_dayofweek(humanized)"] = "Monday"
train_tr.time.loc[train_tr.time["train_tr_dayofweek"] == 1, "train_tr_dayofweek(humanized)"] = "Tuesday"
train_tr.time.loc[train_tr.time["train_tr_dayofweek"] == 2, "train_tr_dayofweek(humanized)"] = "Wednesday"
train_tr.time.loc[train_tr.time["train_tr_dayofweek"] == 3, "train_tr_dayofweek(humanized)"] = "Thursday"
train_tr.time.loc[train_tr.time["train_tr_dayofweek"] == 4, "train_tr_dayofweek(humanized)"] = "Friday"
train_tr.time.loc[train_tr.time["train_tr_dayofweek"] == 5, "train_tr_dayofweek(humanized)"] = "Saturday"
train_tr.time.loc[train_tr.time["train_tr_dayofweek"] == 6, "train_tr_dayofweek(humanized)"] = "Sunday"
print(train_tr.time.shape)
train_tr.time[["TransactionDT_cleaned", "train_tr_dayofweek", "train_tr_dayofweek(humanized)"]].head()
(3540, 403)
```

TransactionDT_cleaned	train_tr_dayofweek	train_tr_dayofweek(humanized)
1970-01-02 00:00:00	4	Friday
1970-01-02 00:00:01	4	Friday
1970-01-02 00:01:09	4	Friday
1970-01-02 00:01:39	4	Friday
1970-01-02 00:01:46	4	Friday



4.2 FE결과

TransactionDT

```
5 rows x 394 columns

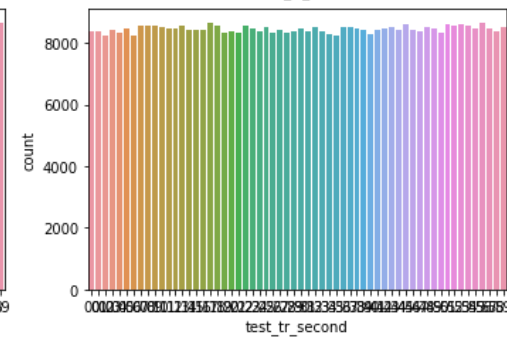
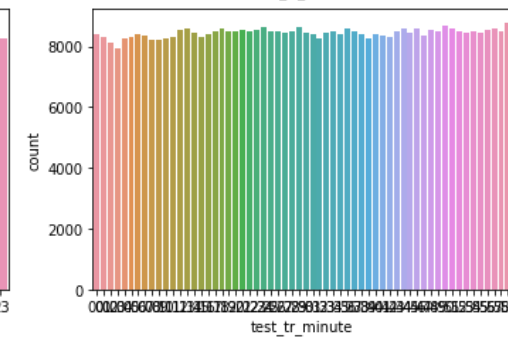
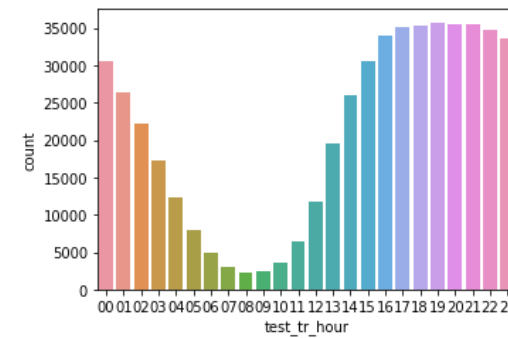
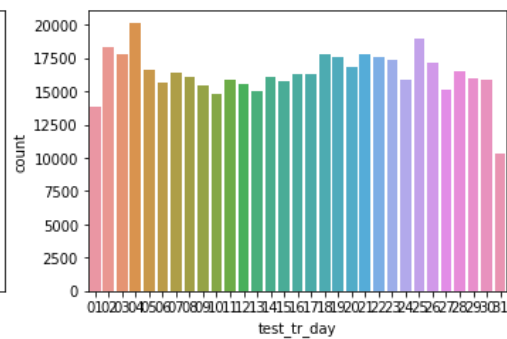
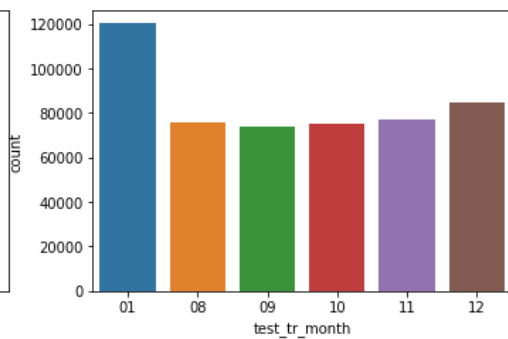
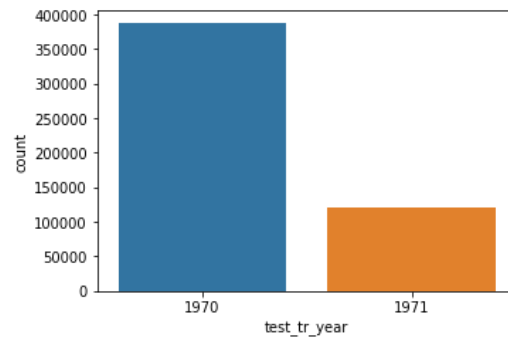
[ ] test_tr_time["test_tr_year"] = test_tr_time["TransactionDT_cleaned"].astype(str).str[:4]
test_tr_time["test_tr_month"] = test_tr_time["TransactionDT_cleaned"].astype(str).str[5:7]
test_tr_time["test_tr_day"] = test_tr_time["TransactionDT_cleaned"].astype(str).str[8:10]
test_tr_time["test_tr_hour"] = test_tr_time["TransactionDT_cleaned"].astype(str).str[11:13]
test_tr_time["test_tr_minute"] = test_tr_time["TransactionDT_cleaned"].astype(str).str[14:16]
test_tr_time["test_tr_second"] = test_tr_time["TransactionDT_cleaned"].astype(str).str[17:]

test_tr_time["test_tr_dayofweek"] = test_tr_time["TransactionDT_cleaned"].dt.dayofweek
print(test_tr_time.shape)
test_tr_time[["test_tr_year", "test_tr_month", "test_tr_day", "test_tr_hour", "test_tr_minute", "test_tr_second", "test_tr_dayofweek"]].head()

(506691, 401)
test_tr_year test_tr_month test_tr_day test_tr_hour test_tr_minute test_tr_second test_tr_dayofweek
0      1970           08         02         00         00         24           6
1      1970           08         02         00         01         03           6
2      1970           08         02         00         01         50           6
3      1970           08         02         00         01         50           6
4      1970           08         02         00         01         57           6

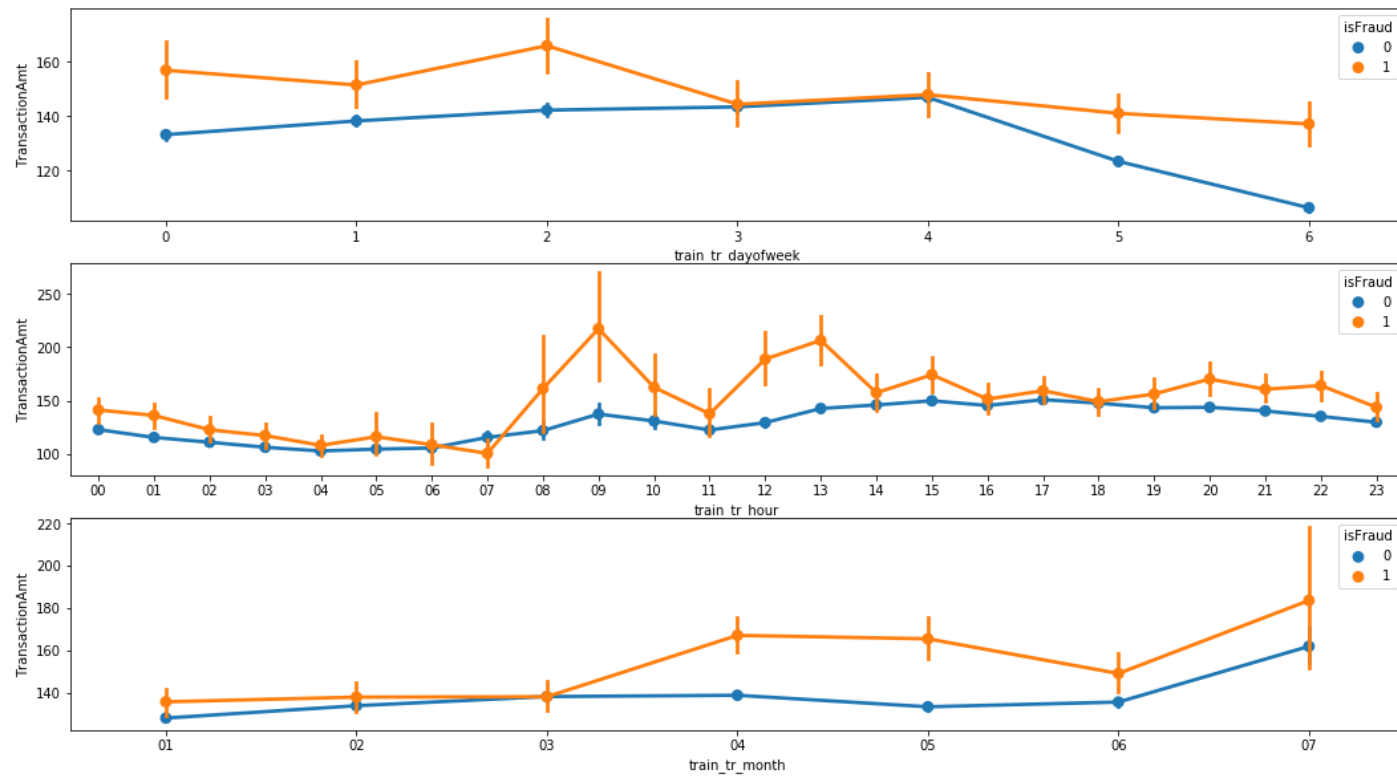
test_tr_time.loc[test_tr_time["test_tr_dayofweek"] == 0, "test_tr_dayofweek(humanized)"] = "Monday"
test_tr_time.loc[test_tr_time["test_tr_dayofweek"] == 1, "test_tr_dayofweek(humanized)"] = "Tuesday"
test_tr_time.loc[test_tr_time["test_tr_dayofweek"] == 2, "test_tr_dayofweek(humanized)"] = "Wednesday"
test_tr_time.loc[test_tr_time["test_tr_dayofweek"] == 3, "test_tr_dayofweek(humanized)"] = "Thursday"
test_tr_time.loc[test_tr_time["test_tr_dayofweek"] == 4, "test_tr_dayofweek(humanized)"] = "Friday"
test_tr_time.loc[test_tr_time["test_tr_dayofweek"] == 5, "test_tr_dayofweek(humanized)"] = "Saturday"
test_tr_time.loc[test_tr_time["test_tr_dayofweek"] == 6, "test_tr_dayofweek(humanized)"] = "Sunday"
print(test_tr_time.shape)
test_tr_time[["TransactionDT_cleaned", "test_tr_dayofweek", "test_tr_dayofweek(humanized)"]].head()

(506691, 402)
TransactionDT_cleaned test_tr_dayofweek test_tr_dayofweek(humanized)
0      1970-08-02 00:00:24           6           Sunday
1      1970-08-02 00:01:03           6           Sunday
2      1970-08-02 00:01:50           6           Sunday
3      1970-08-02 00:01:50           6           Sunday
```



4.2 FE결과

TransactionDT



변형한 TransactionDT 변수들 중
유의미한 fraud차이가 있는 변수는
day of week, hour, month



4.2 FE결과

addr2 & P_emaildomain

Name	TRUE	FALSE	Overall	Percent	Rank	Service info	Region	Country_cod
aim.com	275	40	315	12.6984127	4	AOL	global	1
anonymous.com	36139	859	36998	2.32174712	20	익명	anonymous	9
aol.com	27672	617	28289	2.18105978	25	AOL	global	1
att.net	4003	30	4033	0.74386313	37	미국 통신사	usa	2
bellsouth.net	1856	53	1909	2.77632268	17	미국 통신사	usa	2
cableone.net	156	3	159	1.88679245	29	미국 통신사	usa	2
centurylink.net	205	0	205	0	43	미국 통신사	usa	2
cfl.rr.com	172	0	172	0	43	미국 통신사	usa	2
charter.net	791	25	816	3.06372549	15	미국 통신사	usa	2
comcast.net	7642	246	7888	3.11866126	14	미국 통신사	usa	2
cox.net	1364	29	1393	2.08183776	28	미국 통신사	usa	2
earthlink.net	503	11	514	2.14007782	26	미국 통신사	usa	2
embarqmail.com	251	9	260	3.46153846	11	미국 통신사	usa	2
frontier.com	272	8	280	2.85714286	16	미국 통신사	usa	2
frontiernet.net	190	5	195	2.56410256	19	미국 통신사	usa	2
gmail	485	11	496	2.21774194	23	구글	global	1
gmail.com	218412	9943	228355	4.35418537	9	구글	global	1
gmx.de	149	0	149	0	43	독일 메일회사	germany	4
hotmail.co.uk	112	0	112	0	43	마이크로소프트	uk	3
hotmail.com	42854	2396	45250	5.29502762	8	마이크로소프트	global	1
hotmail.de	43	0	43	0	43	마이크로소프트	germany	4
hotmail.es	285	20	305	6.55737705	6	마이크로소프트	spain	6
hotmail.fr	295	0	295	0	43	마이크로소프트	france	5
icloud.com	6070	197	6267	3.14344982	13	애플	global	1
juno.com	316	6	322	1.86335404	30	미국 통신사	usa	2
live.com	2957	84	3041	2.76224926	18	마이크로소프트	global	1
live.com.mx	708	41	749	5.47396529	7	마이크로소프트	mexico	7
live.fr	56	0	56	0	43	마이크로소프트	france	5
mac.com	422	14	436	3.21100917	12	애플	global	1
mail.com	453	106	559	18.9624329	2	독일 메일회사	germany	4
me.com	1495	27	1522	1.7739816	31	애플	global	1
msn.com	4002	90	4092	2.19941349	24	마이크로소프트	global	1
netzero.com	230	0	230	0	43	미국 통신사	usa	2
netzero.net	195	1	196	0.51020408	39	미국 통신사	usa	2
optonline.net	994	17	1011	1.68150346	32	미국 통신사	usa	2
outlook.com	4614	482	5096	9.45839874	5	마이크로소프트	global	1
outlook.es	381	57	438	13.0136986	3	스페인	spain	6
prodigy.net.mx	206	1	207	0.48309179	40	클라우드펀딩 회사	mexico	7

P_Emaildomain에서
얻어낸
확실한 국가데이터값



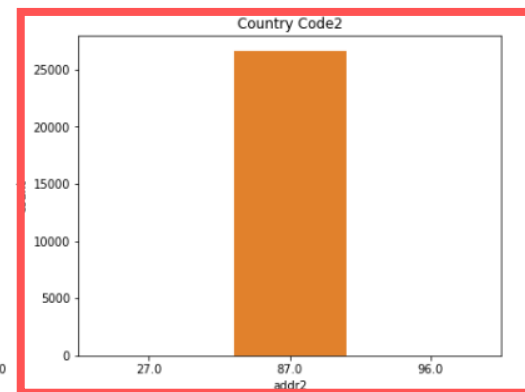
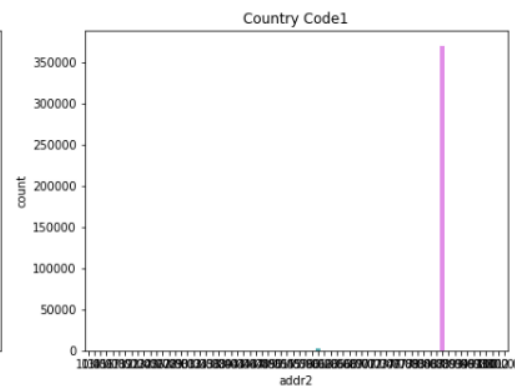
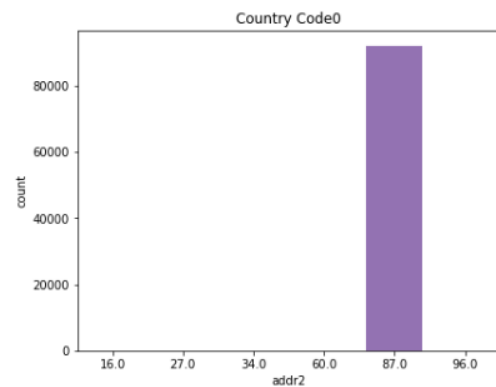
국가코드라
추측되는
Addr2값



4.2 FE결과

addr2 & P_emaildomain

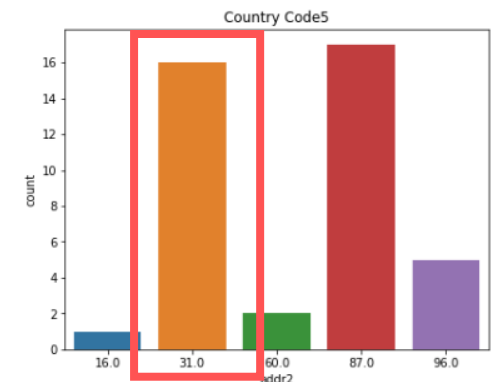
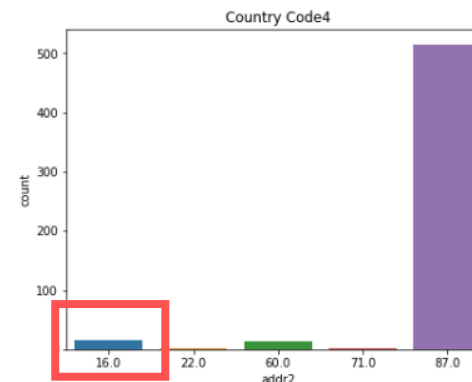
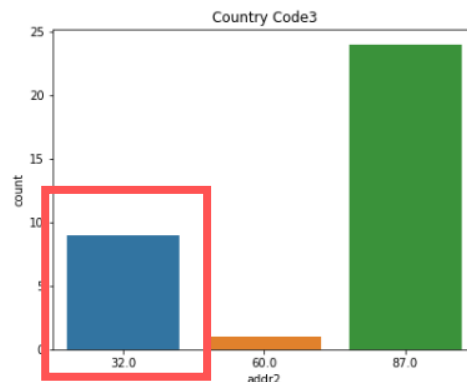
Region from email	Code
Null	0
Global	1
USA	2
UK	3
Germany	4
France	5
Spain	6
Mexico	7
Japan	8
anonymous	9



4.2 FE결과

addr2 & P_emaildomain

Region from email	Code
Null	0
Global	1
USA	2
UK	3
Germany	4
France	5
Spain	6
Mexico	7
Japan	8
anonymous	9



4.2 FE결과

DeviceInfo

train_identity.csv

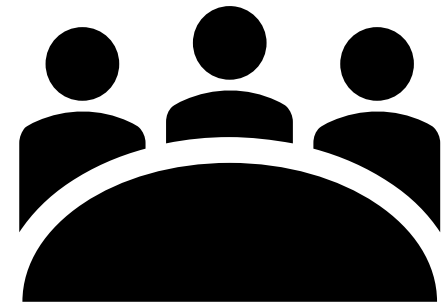
DeviceType

DeviceInfo

id1 - id38

Transaction & Identity

데이터의 통합이 필요함



4.2 FE결과

DeviceInfo

Article Talk Read Edit View history Search Wikipedia

Cramér's V

From Wikipedia, the free encyclopedia

In statistics, **Cramér's V** (sometimes referred to as **Cramér's phi** and denoted as φ_c) is a measure of [association](#) between two [nominal variables](#), giving a value between 0 and +1 (inclusive). It is based on [Pearson's chi-squared statistic](#) and was published by [Harald Cramér](#) in 1946.^[1]

Contents [hide]

- 1 Usage and interpretation
- 2 Calculation
- 3 Bias correction
- 4 See also
- 5 References
- 6 External links

Usage and interpretation [edit]

φ_c is the intercorrelation of two discrete variables^[2] and may be used with variables having two or more levels. φ_c is a symmetrical measure, it does not matter which variable we place in the columns and which in the rows. Also, the order of rows/columns doesn't matter, so φ_c may be used with nominal data types or higher (notably ordered or numerical).

Cramér's V may also be applied to [goodness of fit](#) chi-squared models when there is a $1 \times k$ table (in this case $r = 1$). In this case k is taken as the number of optional outcomes and it functions as a measure of tendency towards a single outcome.^[*citation needed*]

Cramér's V varies from 0 (corresponding to [no association](#) between the variables) to 1 (complete association) and can reach 1 only when the two variables are equal to each other.

isFraud(이분형)데이터와
train_identity.csv 변수들(범주형)
연관성 파악



Cramer's V?

Calculation [\[edit\]](#)

Let a sample of size n of the simultaneously distributed variables A and B for $i = 1, \dots, r; j = 1, \dots, k$

n_{ij} = number of times the values (A_i, B_j) were observed.

The chi-squared statistic then is:

$$\chi^2 = \sum_{i,j} \frac{(n_{ij} - \frac{n_{i.}n_{.j}}{n})^2}{\frac{n_{i.}n_{.j}}{n}}$$

Cramér's V is computed by taking the square root of the chi-squared statistic divided by the sample size :

$$V = \sqrt{\frac{\varphi^2}{\min(k-1, r-1)}} = \sqrt{\frac{\chi^2/n}{\min(k-1, r-1)}}$$

where:

- φ is the [phi coefficient](#).
- χ^2 is derived from [Pearson's chi-squared test](#)
- n is the grand total of observations and
- k being the number of columns.
- r being the number of rows.

The [p-value](#) for the [significance](#) of V is the same one that is calculated using the [Pearson's chi-squared test](#)

The formula for the variance of $V = \varphi_c$ is known.^[3]

In R, the function `cramerV()` from the package `rcompanion`^[4] calculates V using the `chisq.test` function
`cramersV()` from the `lsr`^[5] package, `cramerV()` also offers an option to correct for bias. It applies the

```
[ ] def cramers_v(x, y):
    confusion_matrix = pd.crosstab(x,y)
    chi2 = ss.chi2_contingency(confusion_matrix)[0]
    n = confusion_matrix.sum().sum()
    phi2 = chi2/n
    r,k = confusion_matrix.shape
    phi2corr = max(0, phi2-((k-1)*(r-1))/(n-1))
    rcorr = r-((r-1)**2)/(n-1)
    kcorr = k-((k-1)**2)/(n-1)
    return np.sqrt(phi2corr/min((kcorr-1),(rcorr-1)))
```

```
x = id_new['isFraud']

for i in id_new.columns:
    if i == 'isFraud' or i == 'TransactionID':
        print(i + ' pass')
        continue

y = id_new[i]
```

```
score = cramers_v(x,y)
if score > 0.3:
    print(i + " is significant")
    print(score)
```

```
TransactionID pass
id_21 score:
0.42701968880673
id_25 score:
0.4985777900494445
DeviceInfo score:
0.38344012723979387
isFraud pass
```





1. 목적



2. 변수분석



3. 결측치처리



4. FE



5. 결과



5. 전처리결과

train_transaction.csv



TransactionID TransactionDT

TransactionAmt ProductCD

card1 card2 card3 card4 card5 card6

Addr1 addr2 dist1 dist2

P_emaildomain R_emaildomain

C1-C14 D1-D15 M1-M19 Vxxx



이용변수



TransactionID isFraud
TransactionAmt card1
Card 2 card3 card5

변수변형



TransactionDT DeviceInfo

새로운변수



Addr2 Region

One-hot Encoding



ProductCD card4 card6

기

문승현

11
E

박현지

고예희

염정운