



신용카드 사기 거래 탐지 AI 구축

김세중 김정환 안서연 양재훈 엄소현 정병주



1. About Topic & Data

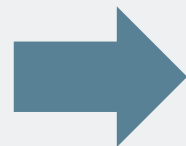
2. Validation Data EDA

3. Baseline Modeling

4. Next Approaches?

About Topic & Data

DACON 신용카드 사기 거래 탐지 AI 경진대회



사기성 신용카드 거래 탐지

해당 경우에 요금 청구 X

사기 거래 사전 탐지

AI 모델 필요



About Topic & Data

DACON 신용카드 사기 거래 탐지 AI 경진대회

Imbalanced

전체 신용카드 거래 건수 대비
사기 거래 건수 매우 적음

Unlabeled

각각에 대한 사기 거래 여부를
정확히 확인하는 건 비효율적



About Topic & Data

DACON 신용카드 사기 거래 탐지 AI 경진대회

Imbalanced

전체 신용카드 거래 건수 대비
사기 거래 건수 매우 적음

Unlabeled

각각에 대한 사기 거래 여부를
정확히 확인하는 건 비효율적

Imbalanced Data에 대한

Unsupervised Binary Classification!

About Topic & Data

평가 기준: Macro-F1 Score

클래스/레이블 별 F1 Score의 평균

➡ 모든 class에 동등한 중요성 부여

➡ Imbalanced한 class가 존재할 경우 사용하기 적절한 메트릭



About Topic & Data

Train Data

- 정상, 사기 여부
Unlabeled
- ID, 비식별화된
feature 30개

Validation Data

- 정상, 사기 여부
Labeled
- ID, 비식별화된
feature 30개,
Class

Test Data

- 정상, 사기 여부
Unlabeled
- ID, 비식별화된
feature 30개



Validation Data EDA

Train Data에는 정상/사기 거래 정보 X



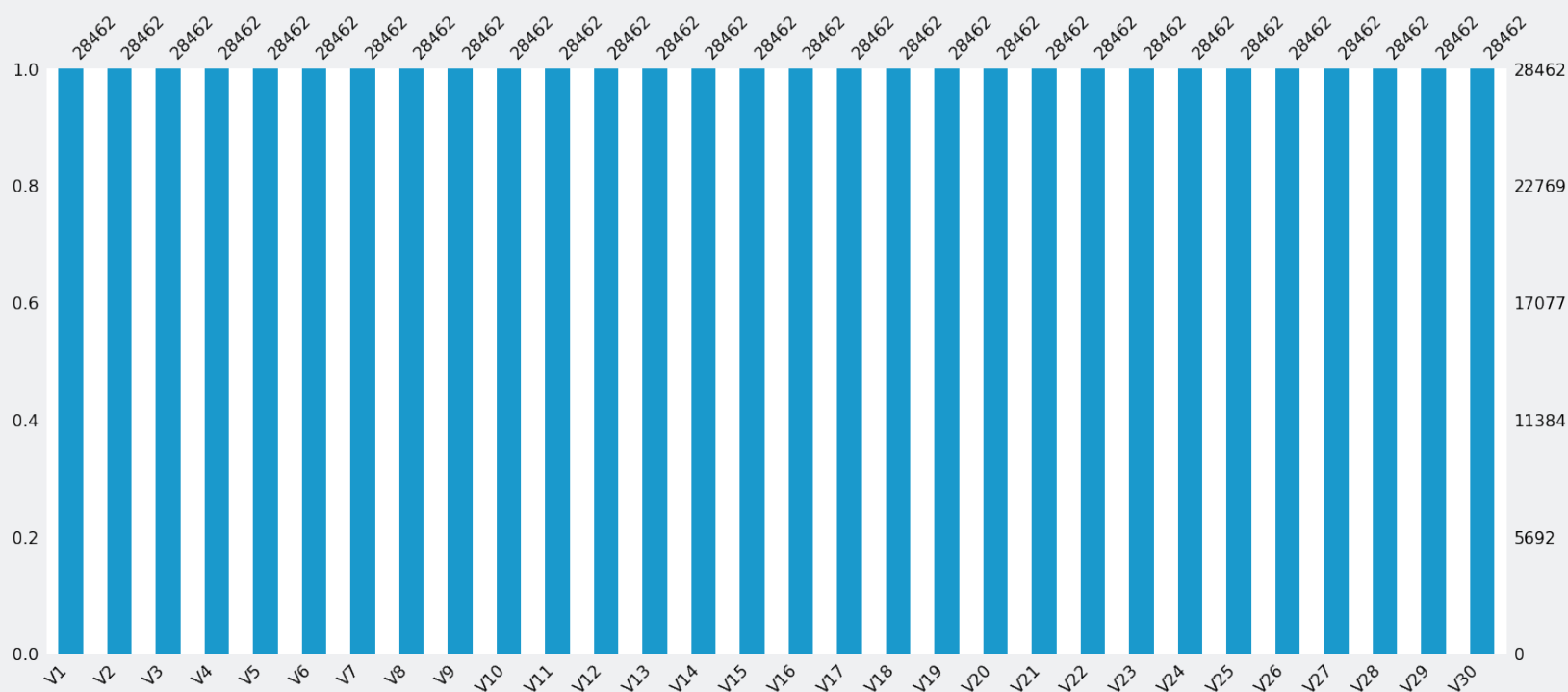
Validation Data EDA를 통해

정상/사기 거래 데이터 간 차이가 있는지 시각적으로 확인해보기 위함!



Validation Data EDA

1. 결측치 여부 확인



결측치가 포함된

데이터 포인트

존재 X

Validation Data EDA

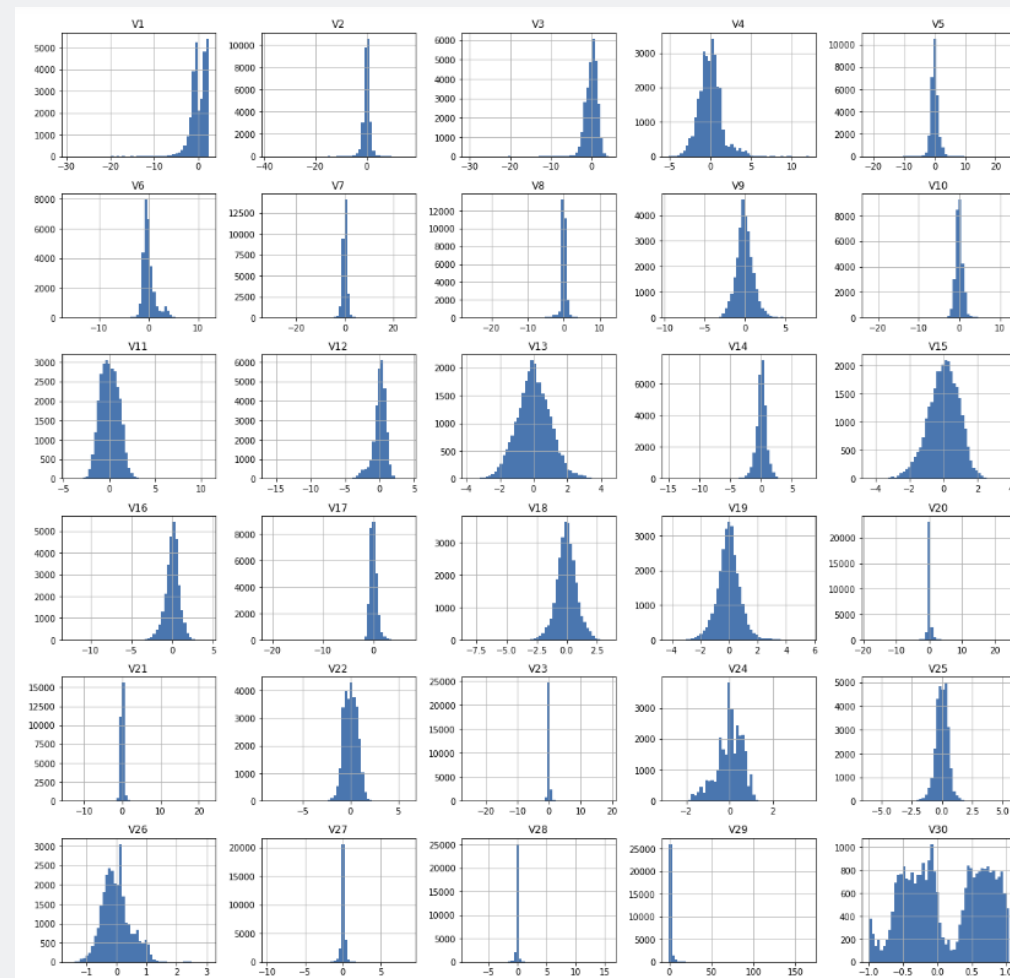
2. Column별 통계량

변수별 히스토그램

각 변수별로 분포 및 값의 폭이 다양함



Scaling 및 정규화 고려



Validation Data EDA

2. Column별 통계량

	V1	V2	V3	V4	V5	V6	V7	V8	...	V22	V23	V24	V25	V26	V27	V28	V29	V30
count	28432	28432	28432	28432	28432	28432	28432	28432	...	28432	28432	28432	28432	28432	28432	28432	28432	28432
mean	0.01182	-0.00303	0.011758	-0.00387	0.0009	-0.01581	0.013685	0.007648	...	0.000826	0.001394	0.001554	0.002346	0.000435	0.000032	-0.000086	0.924202	0.120106
std	1.903021	1.592242	1.438086	1.390272	1.311665	1.290693	1.101884	1.092214	...	0.720846	0.587799	0.603244	0.526565	0.484707	0.38878	0.30378	3.348744	0.558146
min	-29.5161	-38.3053	-14.8481	-5.07124	-21.577	-16.1726	-16.3871	-26.278	...	-8.55581	-25.3567	-2.8079	-6.03505	-1.59649	-9.79357	-8.364853	-0.30741	-0.99488
0.25	-0.91247	-0.59885	-0.86995	-0.8531	-0.69688	-0.77374	-0.54483	-0.21102	...	-0.54276	-0.1603	-0.35223	-0.32208	-0.32865	-0.07176	-0.052335	-0.22637	-0.35768
0.5	0.025035	0.074632	0.177232	-0.02408	-0.05274	-0.28029	0.047429	0.022312	...	0.008156	-0.01135	0.040176	0.016615	-0.04977	0.000434	0.012136	0.006218	0.002408
0.75	1.315867	0.801268	1.012016	0.735528	0.598752	0.378145	0.56738	0.323227	...	0.527964	0.145999	0.441544	0.353035	0.24052	0.088949	0.080736	0.777754	0.64186
max	2.411769	11.87496	4.226108	9.752791	24.34531	12.12895	26.23772	11.54125	...	6.090514	18.94673	3.658746	5.525093	3.067907	8.708972	15.726807	165.9483	1.034975

변수별로 Min, Max 값의 차이가 매우 크기 때문에 **scaling** 필요성 재확인

>>> Validation Data EDA

2. Column별 통계량

변수별 box plot 통해 이상치 존재 확인

But 이상치 데이터들이

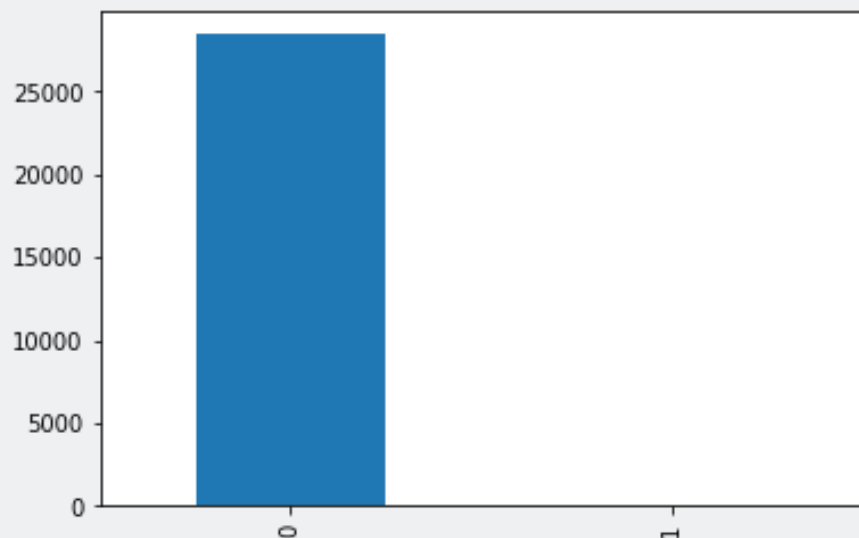
사기 거래 탐지의 key point가 될수도?

일단 제거하지 않기로



Validation Data EDA

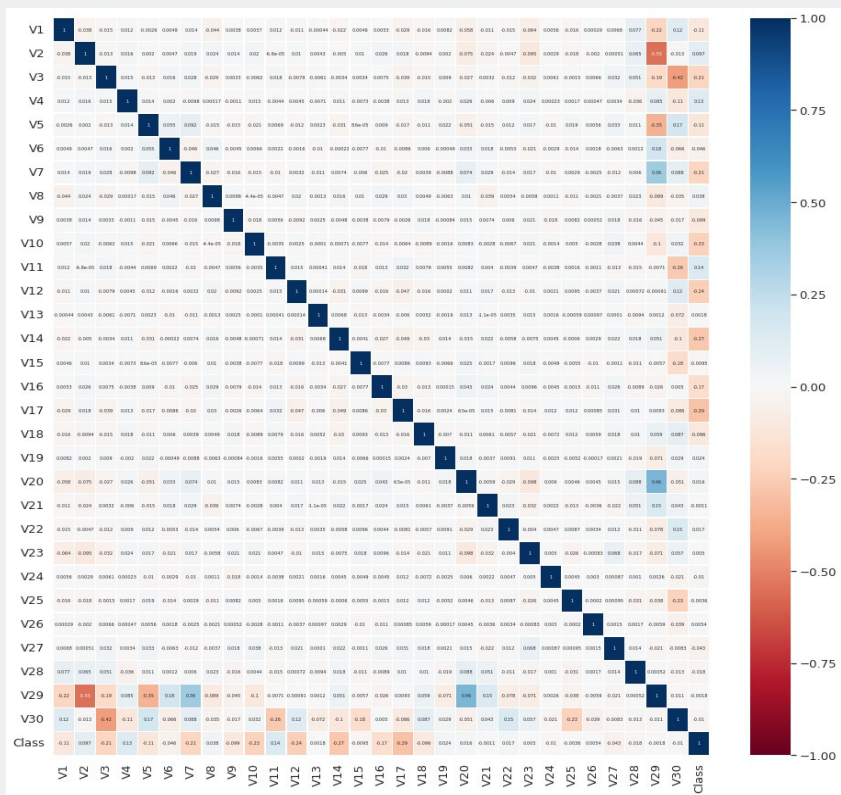
3. Target 변수 분포



Detecting 해야 하는
사기 거래의 건수가 현저히 적음
사기 거래 발생 비율 = 0.11%



4. Heatmap



Validation Data에 대한 Heatmap

V17, V14, V12, V10, V7, V3 변수가

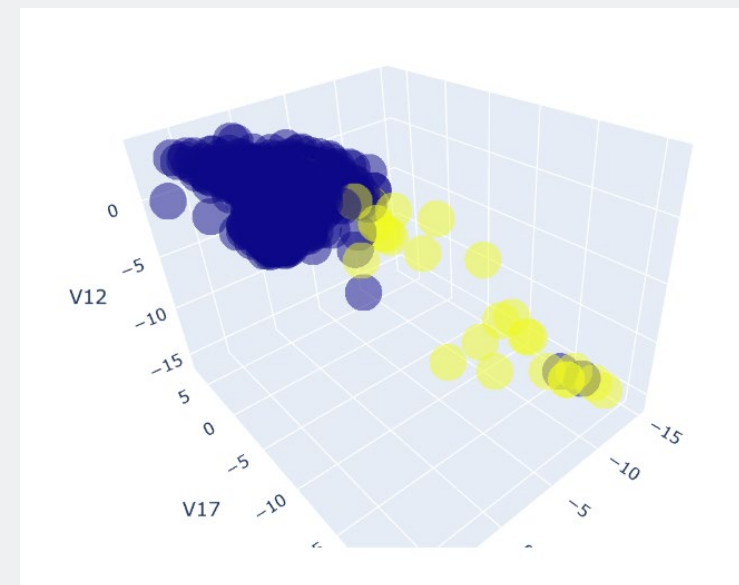
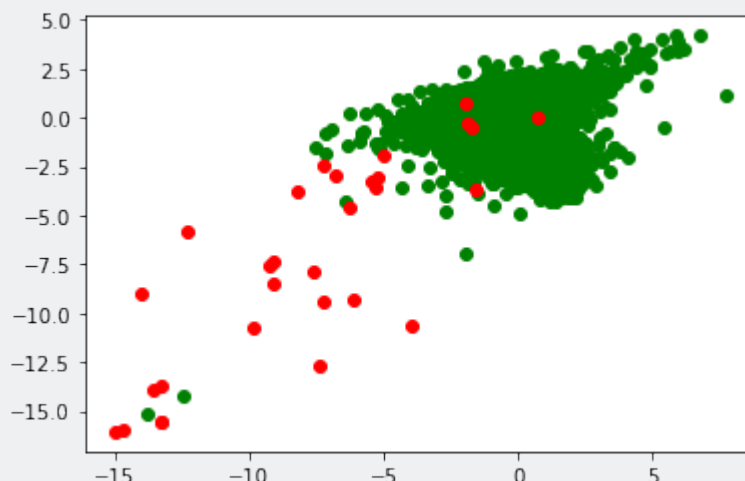
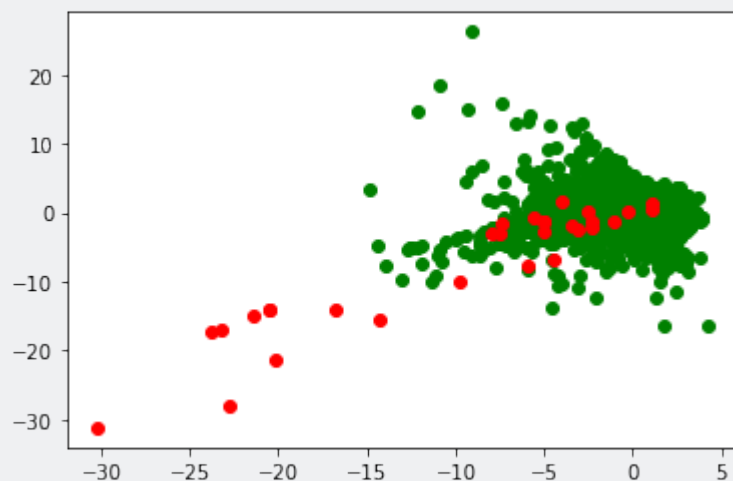
절댓값 0.2 이상으로 Class(정상/비정상)과

매일 correlation을 보임

Feature간 상관관계는 높지 않음

Validation Data EDA

4. Scatter Plot



Correlation이 높았던 변수들 위주로 2, 3차원 scatter plot

정상 데이터끼리 모여 있음, 3개 변수 만으로도 class가 꽤 가시적으로 분류됨

Baseline Modeling

~~Imbalanced Data
Handling~~

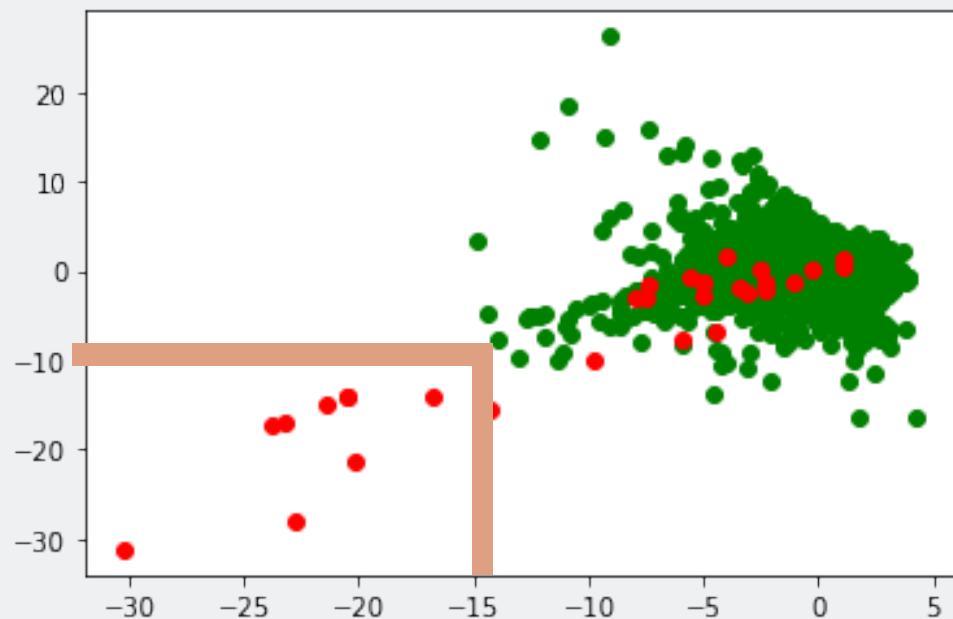
~~Ensemble 등
고도화된 모델링 방법~~

기본적인 단일 model의 성능을 알아보고
이를 baseline으로 이용하려는 의도로 진행



Baseline Modeling

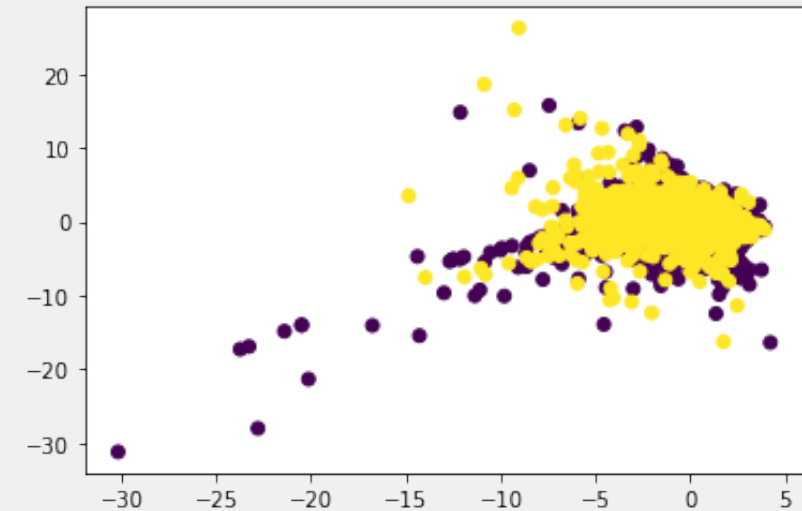
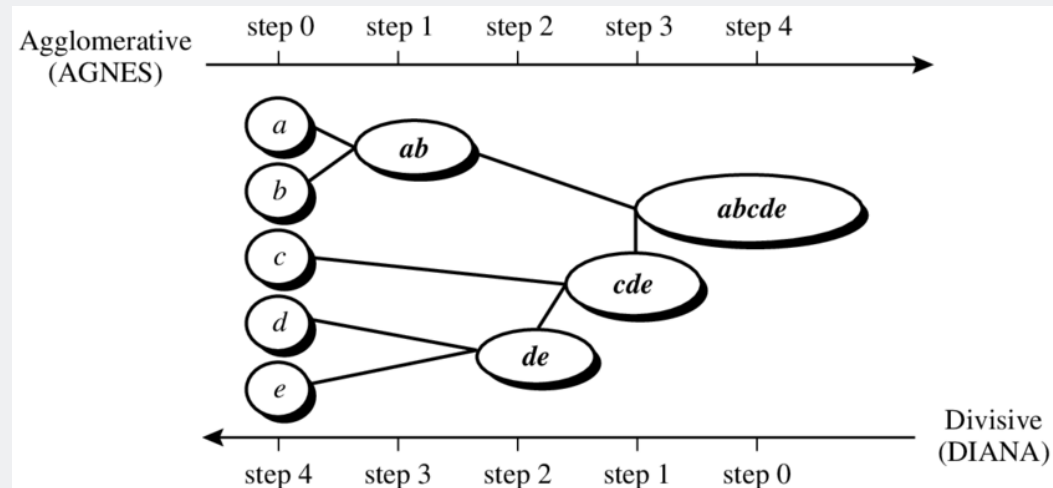
1. EDA를 통한 Simple Conditional Model



Validation Macro f1 score	0.6890
Test(30%) Macro f1 score	0.5684
Test(70%) Macro f1 score	0.6030

Baseline Modeling

2. AGglomerative NESTing(AGNES)

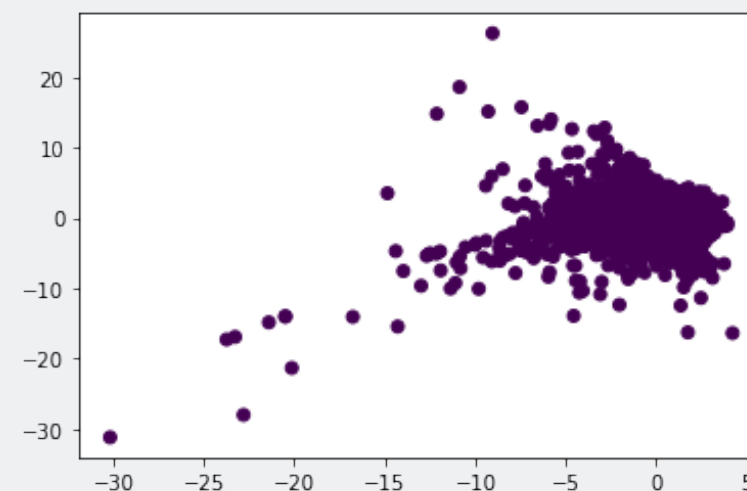
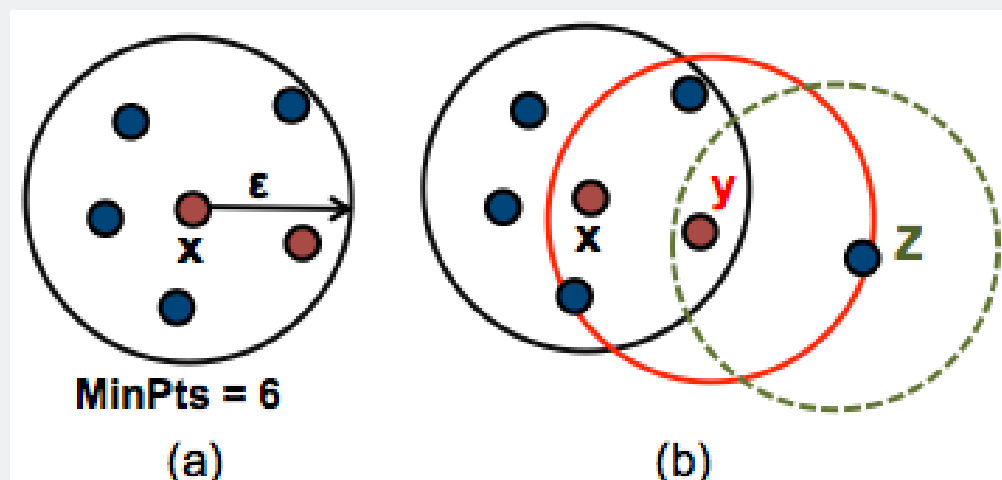


Validation
Macro f1 score

0.3519

Baseline Modeling

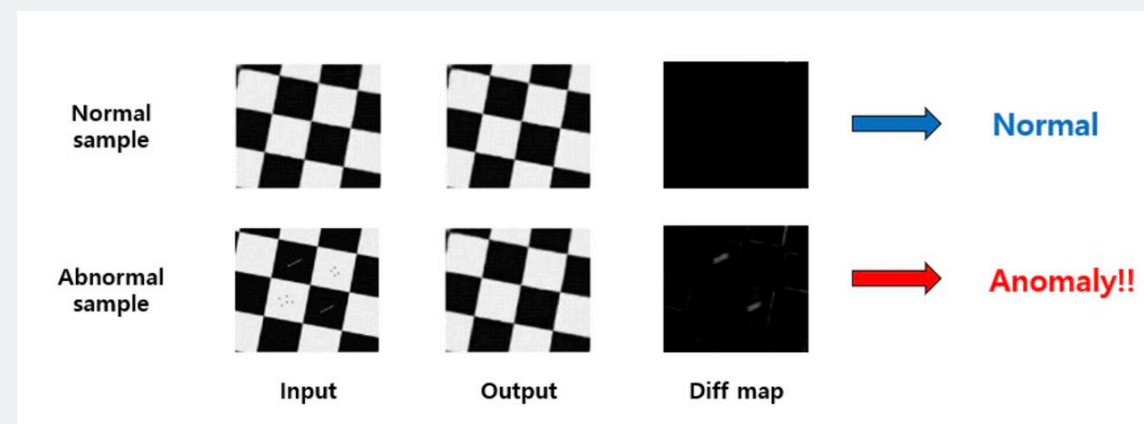
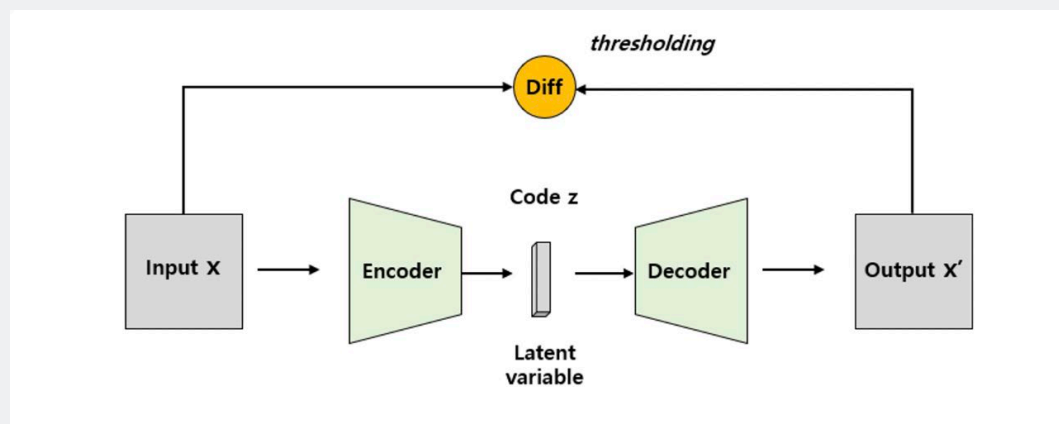
3. DBSCAN



Validation Macro f1 score	0.4997
------------------------------	--------

Baseline Modeling

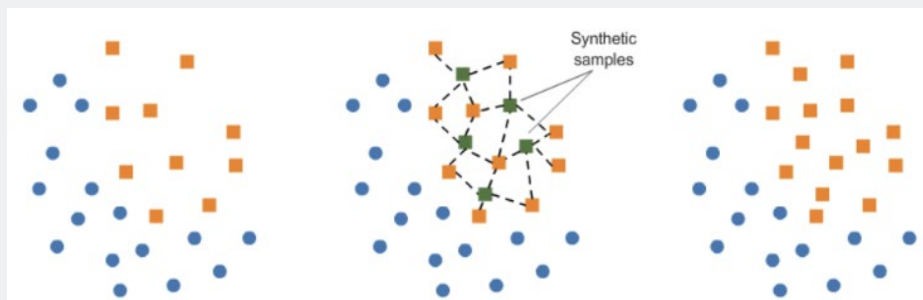
4. Auto Encoder



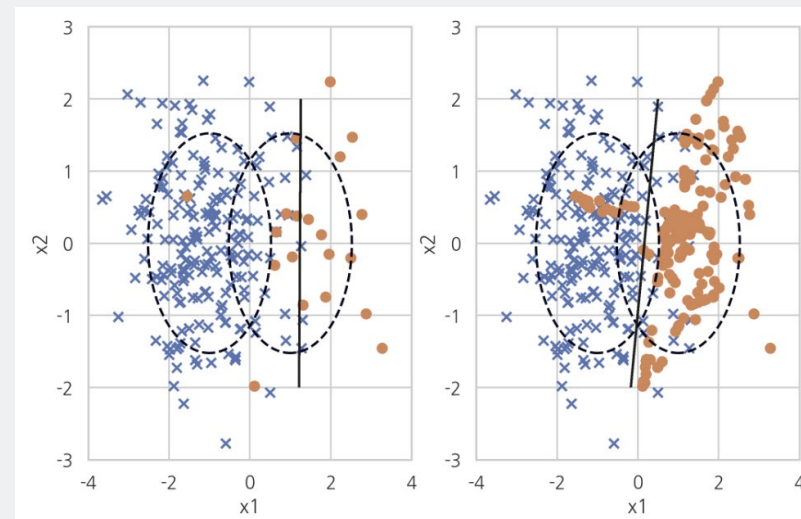
Validation
Macro f1 score

0.9166

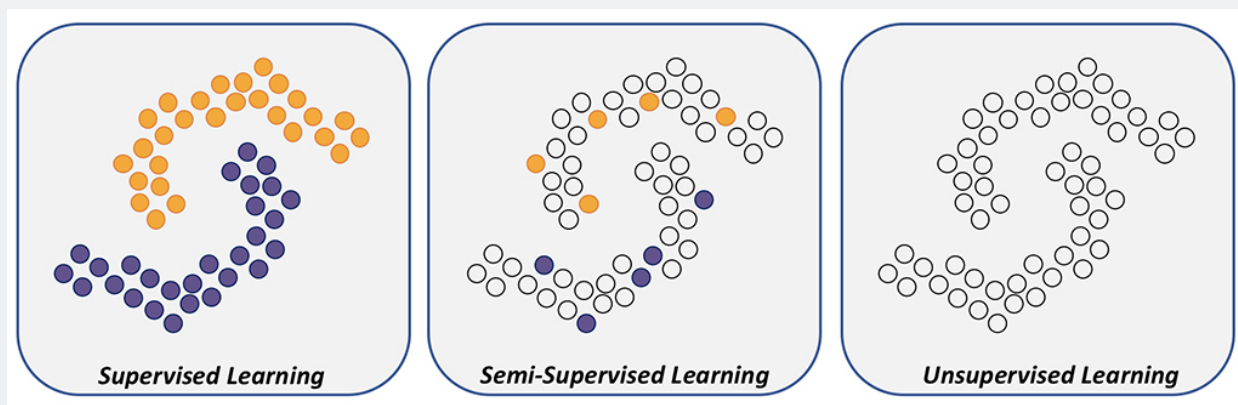
Next Approaches?



Over Sampling-SMOTE



Over Sampling
-Adaptive Synthetic
Sampling



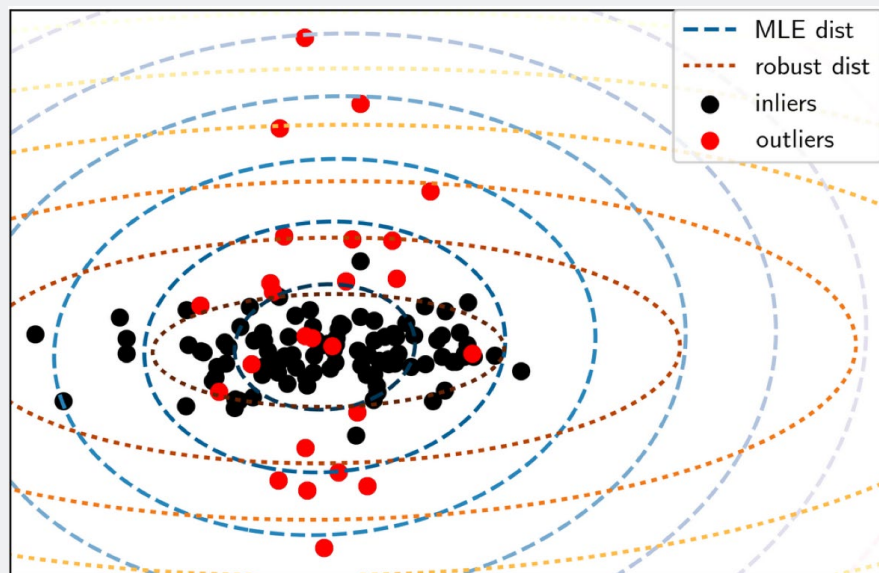
Semi Supervised Learning

Next Approaches?

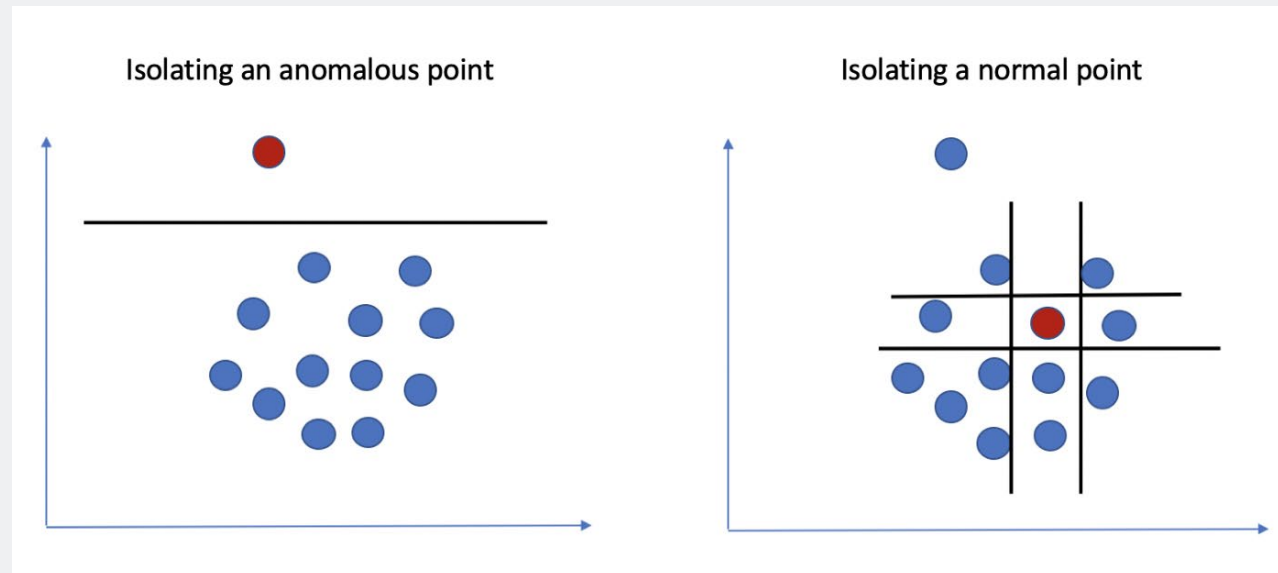


Next Approaches?

Anomaly Detection



EllipticEnvelope



Isolation Forest



End of Presentation
