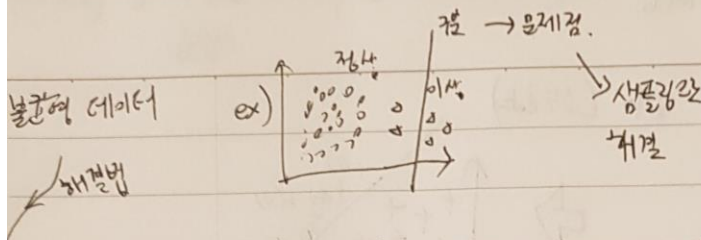


불균형 data 처리

No. 20220307 고종국



샘플링 기법

언더샘플링: 샘플개수를 축소.
!계산시간↓ !모델오버피팅↓
!정보손실↑

오버샘플링: 샘플개수 증가
!정확도↑ !모델오버피팅↓

모델링 기법

!과적합가능성↑
!계산시간↑
!노이즈 민감

Random undersampling

Tomek link

Condensed Nearest Neighbor Rule

One-sided selection

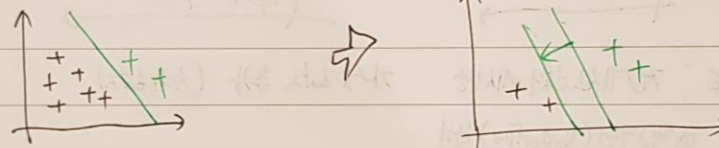
Resampling

SMOTE

Borderline-SMOTE

ADASYN

Random Sampling: 다수 data 에서 랜덤으로 샘플링



Tomek link: 다수 data 에서 Tomek link 형성된 것 (경계선에 있는 data) 제거

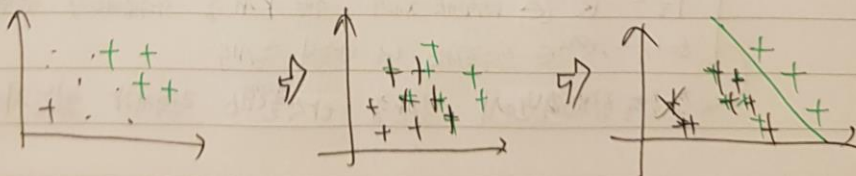


Condensed nearest neighbor

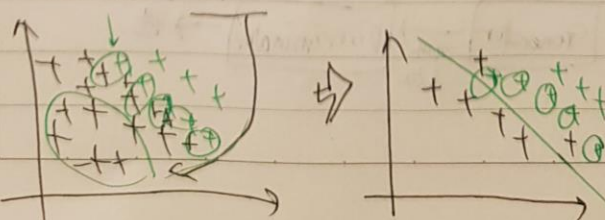
Neighbor
⇒ CNN

① 다수 data 에서 개만 남김
② 1/NN경리가 원래 다수 data 로
직접것들 다 제거

③ 나머지 다수 data 각각의 위치에서
1/NN



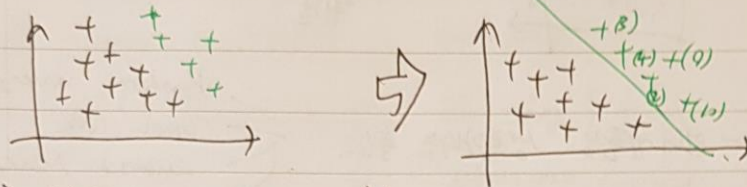
One Side Selection: Tomek link + CNN



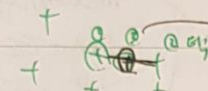
4) 오차율 기법으로
No. 불균형 문제 해결

비율기반학습 : 소수 data에 가중치 부여.
단일 클래스 분기법 : 보통 다수 data를, 그걸의 분기점계선 (ex. 원형) 생성

Resampling : 같은 데이터 중첩해서 증가 (소수 data)

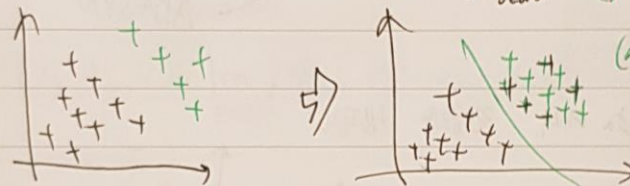


SMOTE : 소수 data에서 가상 data 추가

if $k=5$,  \rightarrow ① ② 방향으로 가상 data ③ 추가.

$$\text{가상 data 생성} = \underbrace{x_i}_{\text{①}} + \underbrace{u}_{\text{②}} \cdot (\underbrace{x_{(n)}}_{\text{③}} - x_i)$$

(0~1) 랜덤. ① ② ③

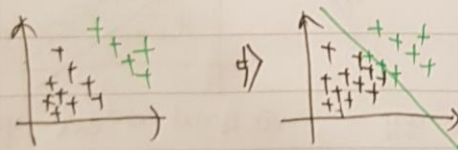


Borderline-SMOTE : 경계선 근처에서만 가상 data 추가 (소수 data)

① 경계선 (Borderline) 찾기

② 소수 클래스 중 k 는 가장, 그중 k 개 주변 탐색

③ k 개 중 다수 클래스 수 k' 기준, 1) $k < k'$ \rightarrow Noise. ② $k < k'$ \rightarrow Borderline OK. 3) $k > k'$ \rightarrow Safe 할 것.



ADASYN : 위치마다 다르게 가상 data 추가 (소수 data)

$$r_i = \frac{\Delta D_i}{K}$$
 \rightarrow 소수 클래스 x_i 의 주변 K 개 중 다수 클래스의 관측치 개수
 $i=1, \dots, m \leftarrow$ 소수 클래스 내 관측치 총 개수
 \rightarrow r_i 소수 클래스 주변에 얼마나 많은 다수 클래스 관측치가 있는지에 대한 지표.

4) GAN 응용

PUA 생성 알고리즘

