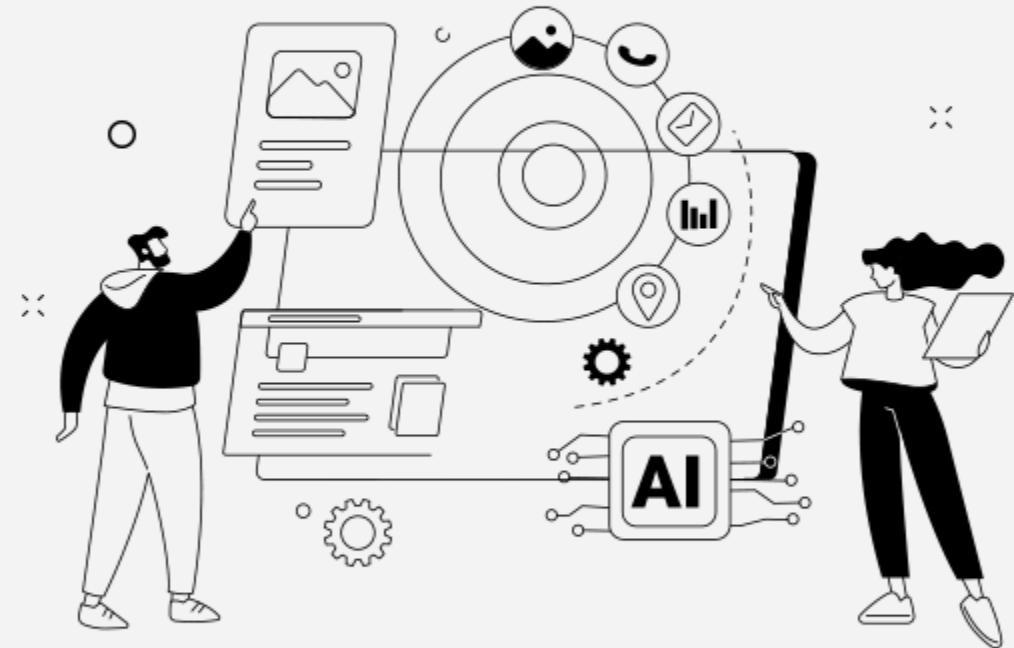


2022 데이터 크리에이터 캠프

Data Creator Camp



대학팀 0824



과학기술정보통신부

NIA 한국지능정보사회진흥원

팀 소개



연준

EDA



재우

실사 이미지 제거



세연

데이터 불균형



윤희

모델 구축

팀 협업 방법

Data Creator Camp - Team 0824

프로젝트 차트

A 이름	B 담당자	C 주자
미션1 - EDA 시작회	김연준	1주차
미션3 - 코랩 파일 정리	Seyeon Park	박윤희 박재우 김연준
제출파일에 모델 테스트		5주차
PPT 제작	Seyeon Park	박윤희 박재우 김연준
미션3 - 모델링(Resnet, 이미지 모델 실험)	박윤희	4주차
미션3 - 모델링, 성능평가(최종)	박재우 박윤희	5주차
미션2 - 코랩 파일 정리	Seyeon Park	박윤희 박재우 김연준
미션3 - 모델링(기본)	Seyeon Park	박윤희 박재우 김연준
미션2 - 실사 이미지 제거(DBSCAN)	박재우	3주차
미션2 - 실사 이미지 제거(DBSCAN) Filtering	김연준 박재우	3주차
미션2 - 실사 이미지 제거(K-means)	박재우	3주차
미션2 - 실사 이미지 제거(Object Detection)	Seyeon Park	3주차
미션1 - EDA	Seyeon Park	박윤희 박재우 김연준
미션2 - 코랩 파일 정리	Seyeon Park	박윤희 박재우 김연준
미션2 - imbalanced data 해결	Seyeon Park	2주차

미팅 4

2022년 10월 14일

⑤ Seyeon Park ⑨ 박윤희 ⑨ 박재우 ⑨ 김연준

미션2 방법 논의, 엔도린 준비

<1014_금_미팅4>

- ▶ 리얼 이미지 특징 정리
- ▶ 클러스터링 기법
- ▶ 중복 이미지
- ▶ 과정도 포함해 결과 제출하자는 의견

<멘토링 질문>

- EDA 어느 정도까지 해야 하는지?
- 중복 이미지 제거에 대한 조언 → 무조건 제거해야 하나? (필요성)
- 무조건 중복 이미지를 제거해야 하는지?
- 리얼 이미지 제거에 대한 조언
- 이미지 크기 차이에 대한 조언
- 만들었던 AI를 전이 학습 시에만 이용하지 못하는 건지, 전체적으로 이용 못하는 건지?
- AI기술이 어떤 범위까지를 말하는 건지?
- 결과 제출 양식은 어떻게 되는지?

회의록

미팅 8 2022년
미팅 9 2022년

Google Drive

드라이브

내 드라이브 > 2022_DataCreator - Google Sheets

모델결과정리

	A	B	C	hyperparameters	D
	data	model			acc
1	Random Undersampling + 48% Upsampling	resnet	clipX, AdamW, lr=5e-5		0.691
2	Random Undersampling + 48% Upsampling	resnet	clipX, Adam, lr=1e-3		0.611
3	Random Undersampling + 48% Upsampling	resnet	clipO, Adam, lr=1e-3		0.661
4	Random Undersampling + 48% Upsampling	resnet	clipO, Adam, lr=2e-4		0.695
5	Random Undersampling + 48% Upsampling	resnet	clipO, AdamW, lr=2e-4		0.611
6	Random Undersampling + 48% Upsampling	resnet	clipO, AdamW, lr=5e-5, augmentX		0.605
7	Random Undersampling + 48% Upsampling	baseline model	Adam, lr=1e-6, dropout 0.2 + 규제 완화		0.776
8	Random Undersampling + 48% Upsampling	baseline model	Adam, lr=1e-3, reduce_lr		0.855
9	Random Undersampling + 48% Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_O 평가규정*		0.612
10	Random Undersampling + 48% Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_X 평가규정*		0.624
11	Random Undersampling + 48% Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_O 평가규정*		0.647
12	Random Undersampling + 48% Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_X 평가규정*		0.512
13	Random Undersampling + 48% Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_O 평가규정*		0.641
14	Random Undersampling + 48% Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_X 평가규정*		0.512
15					
16	Random Undersampling + 비율대로(2,34) Upsampling	resnet	clipX, AdamW, lr=5e-5		0.691
17	Random Undersampling + 비율대로(2,34) Upsampling	resnet	clipX, Adam, lr=1e-3		0.634
18	Random Undersampling + 비율대로(2,34) Upsampling	resnet	clipO, Adam, lr=1e-3		0.615
19	Random Undersampling + 비율대로(2,34) Upsampling	resnet	clipO, Adam, lr=2e-4		0.595
20	Random Undersampling + 비율대로(2,34) Upsampling	resnet	clipO, AdamW, lr=2e-4		0.641
21	Random Undersampling + 비율대로(2,34) Upsampling	resnet	clipO, AdamW, lr=5e-5, augmentX		0.605
22					
23	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, dropout 0.2 + 규제 완화		0.776
24	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-3, reduce_lr		0.612
25	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_O 평가규정*		0.695
26	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_X 평가규정*		0.612
27	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_O 평가규정*		0.655
28	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_X 평가규정*		0.612
29	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, dropout 0.2 + 규제 완화 + 스케줄러 추가		0.78
30	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, reduce_lr + 스케줄러 추가		0.82
31	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_O 평가규정* 스케줄러 추가		0.747
32	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_X 평가규정* 스케줄러 추가		0.811
33	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_O 평가규정* 스케줄러 추가		0.78
34	Random Undersampling + 비율대로(2,34) Upsampling	baseline model	Adam, lr=1e-4, reduce_lr, augment_X 평가규정* 스케줄러 추가 2		0.841



목차



EDA

- 데이터 셋 살펴보기
- 문제점 도출
- 해결 방안



실사 이미지 제거

- 객체 탐지
- 군집화
- 필터링



데이터 불균형 해결

- Oversampling
- Undersampling
- 클래스 가중치



모델링 및 성능평가

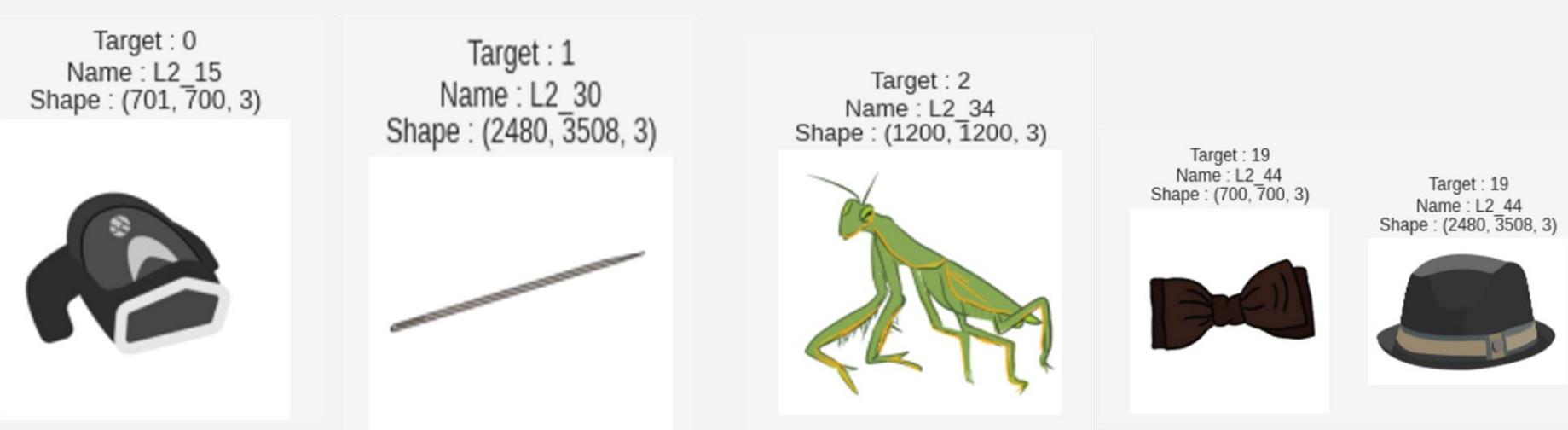
- Resnet 기반 모델
- 커스텀 모델
- 성능평가



과학기술정보통신부

NIA 한국지능정보사회진흥원

데이터 셋 살펴보기



L2_15: 전자기기

L2_30: 생활용품

L2_34: 곤충

L2_52: 액세서리

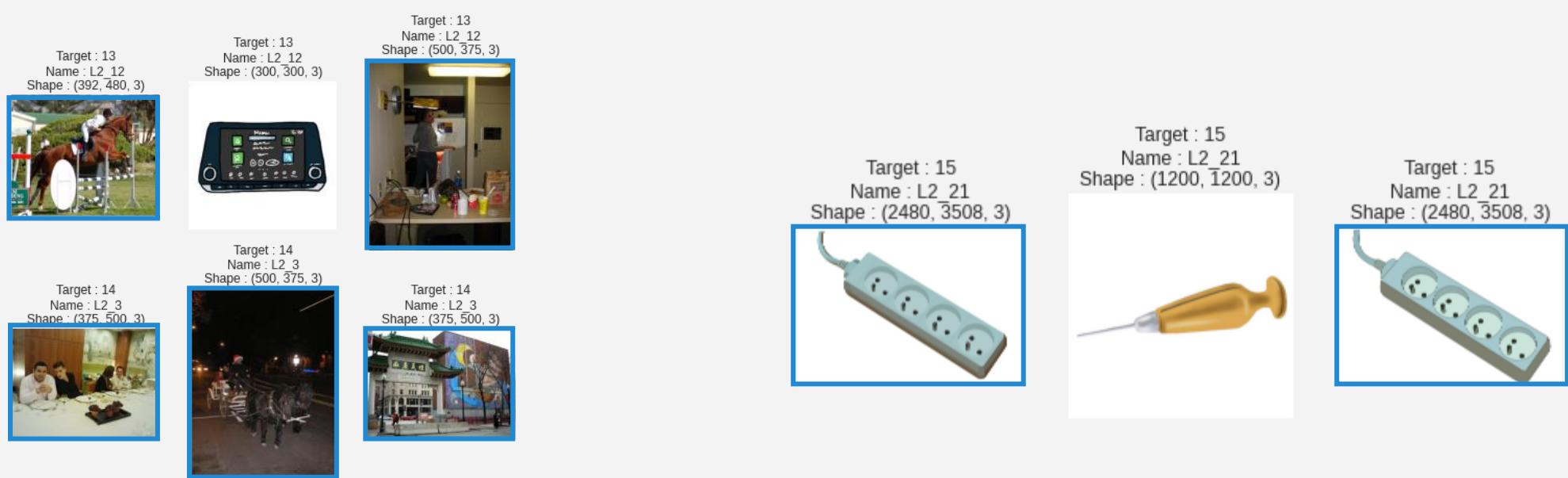
이미지 데이터 25503장, 클래스 20개



과학기술정보통신부

NIA 한국지능정보사회진흥원

문제점 도출

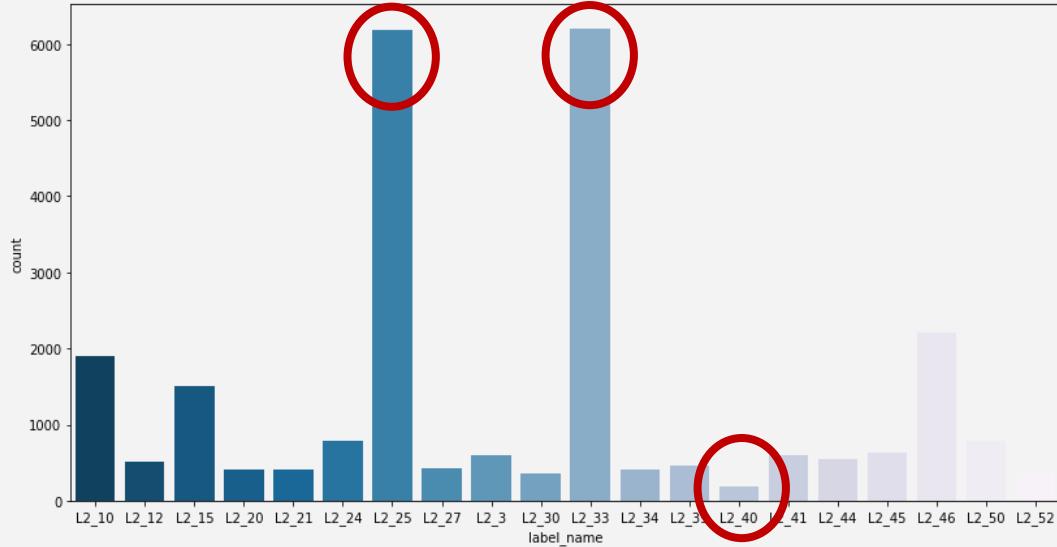


실사 이미지(Real Image)

중복 혹은 유사 이미지

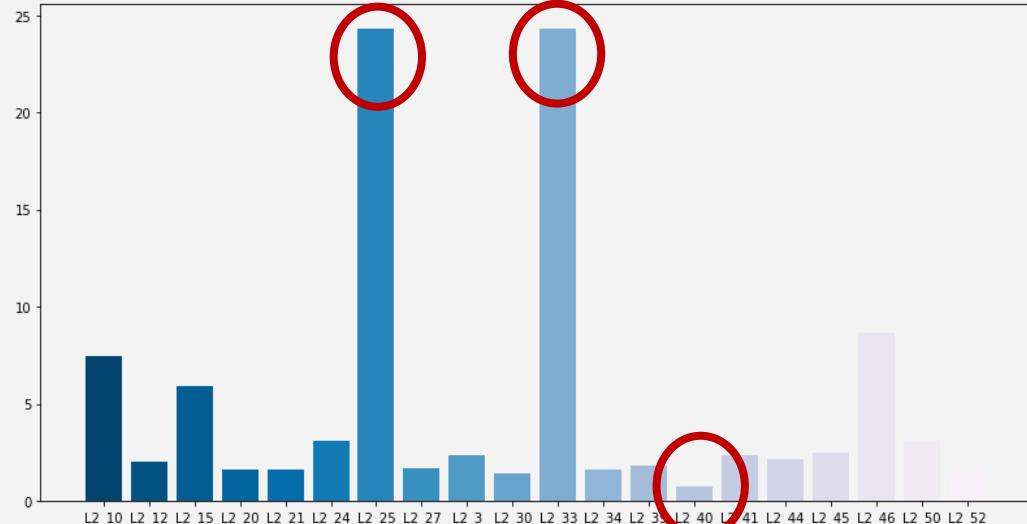
문제점 도출

Distribution by label in dataset



클래스별 데이터 개수 분포

Distribution by label ratio in dataset (%)



클래스별 데이터 비율 분포

데이터 불균형

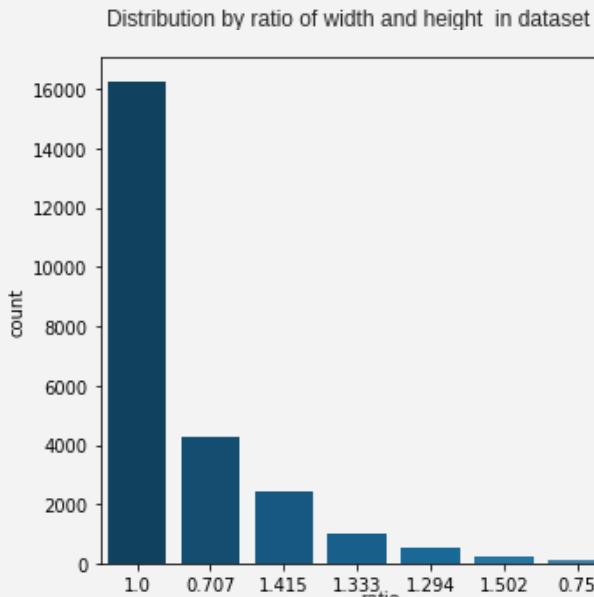
L2_33(6206개)가 클래스 L2_40(180개) 보다 **약 34배** 정도 많음



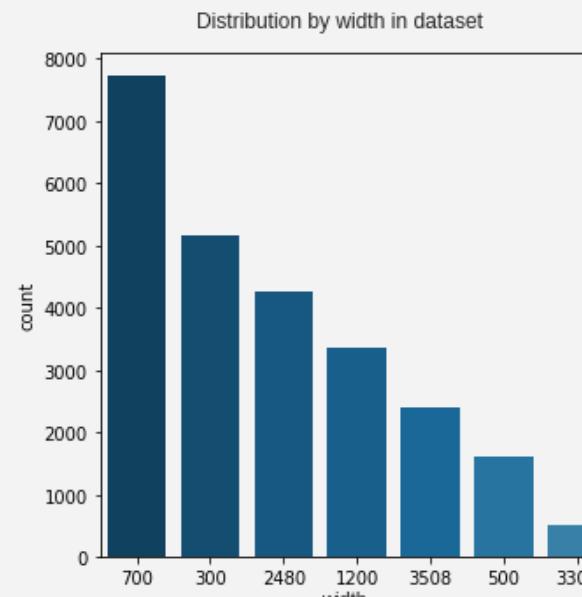
과학기술정보통신부

NIA 한국지능정보사회진흥원

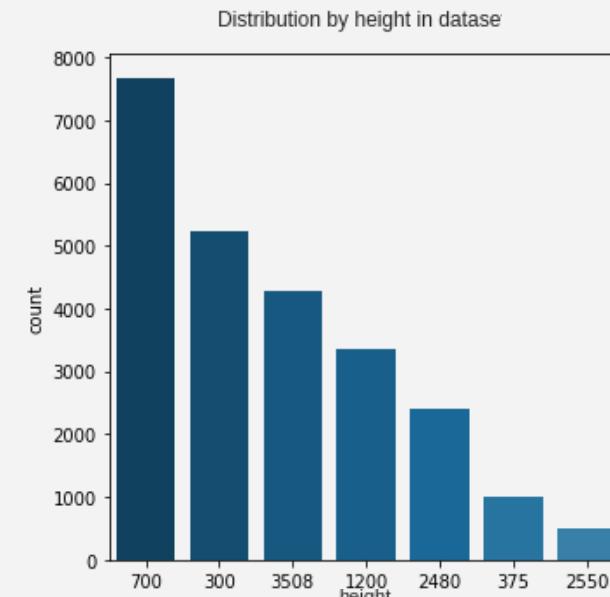
문제점 도출



이미지의 가로 세로 비율 분포



이미지의 가로 사이즈 분포



이미지의 세로 사이즈 분포

이미지별 사이즈 불일치



과학기술정보통신부

NIA 한국지능정보사회진흥원

해결방안

문제점

실사 이미지 존재

중복 이미지 존재

클래스별 데이터 불균형

이미지별 사이즈 불일치



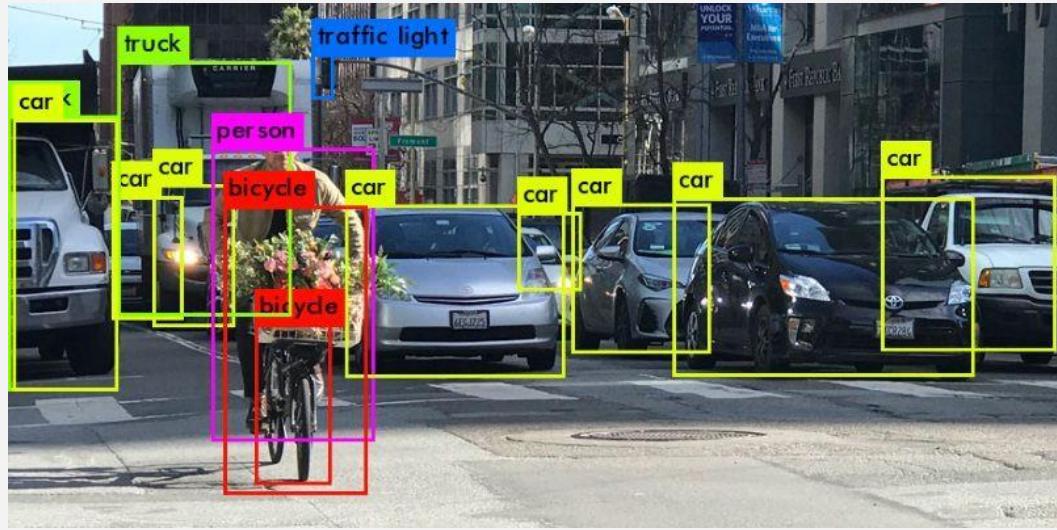
해결방안

- | | | | |
|--------------------|----------|-----------------|------------------------|
| - Object Detection | - 이미지 증강 | - Oversampling | - 파라미터 조절
(256,256) |
| - Clustering | | - Undersampling | |
| - Filtering | | - 클래스 가중치 | |

실사 이미지 제거

Solution 1 - Object Detection

Object Detection 모델 Yolo를 사용해 탐지되는 객체 수에 따라 이미지 제거



Solution 2 - Clustering

이미지의 Feature로 군집화하고 실사 이미지 군집은 제거

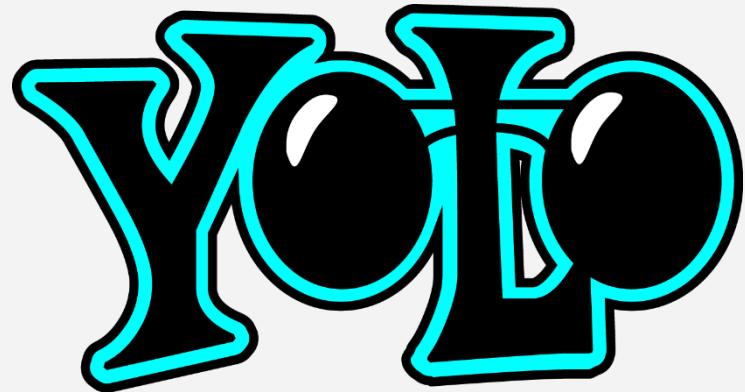


과학기술정보통신부

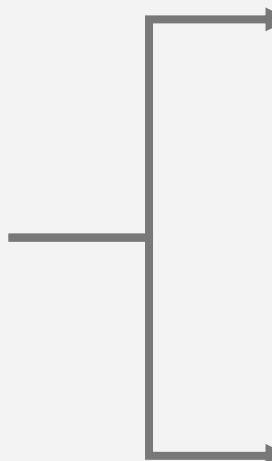
NIA 한국지능정보사회진흥원

객체 탐지

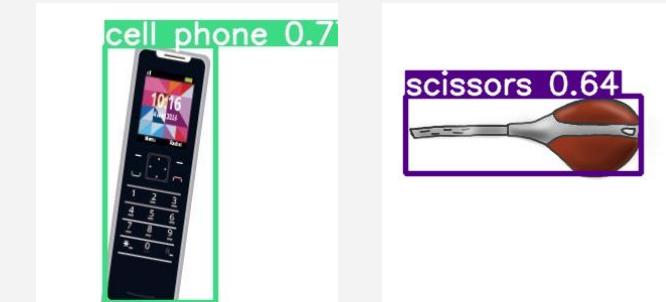
Solution 1 - Object Detection



Object Detection 모델 YOLO로
이미지를 Detect



객체가 1개일 경우: 일러스트



객체가 2개 이상일 경우: 실사 이미지



과학기술정보통신부

NIA 한국지능정보사회진흥원

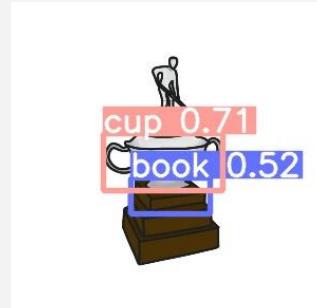
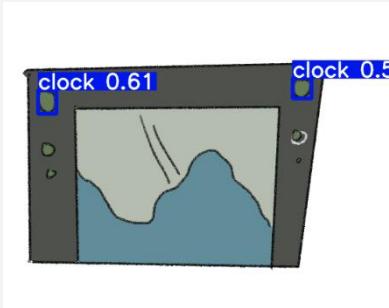
객체 탐지

Solution 1 - Object Detection

객체가 1개지만 실사 이미지



객체가 2개 이상이지만 일러스트



Solution 2 Clustering



일러스트인지? 실사 이미지인지?
라벨이 필요 없는 비지도 학습

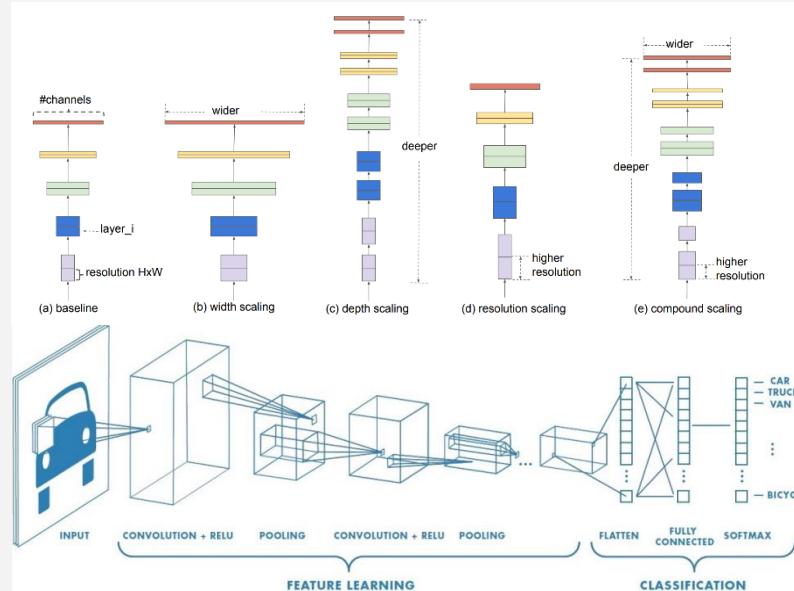


과학기술정보통신부

NIA 한국지능정보사회진흥원

군집화

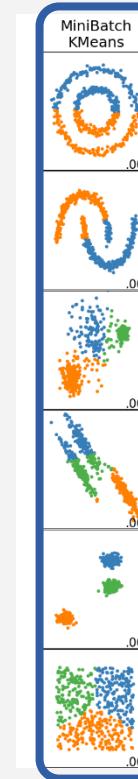
Solution 2 - Clustering



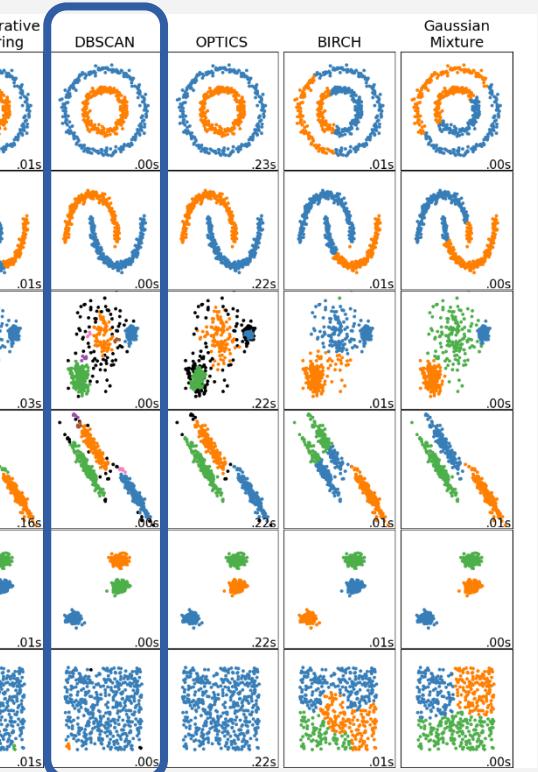
Efficientnet b4

CNN based
Feature Extractor

KMeans



DBSCAN



Clustering model

군집화 모델 비교

Solution 2 - Clustering

Best Clustering Model

KMeans < DBSCAN

WHY?

KMeans

구역을 기준으로
Clustering



DBSCAN

각 원소들의 군집 밀도에 따라
Clustering



과학기술정보통신부

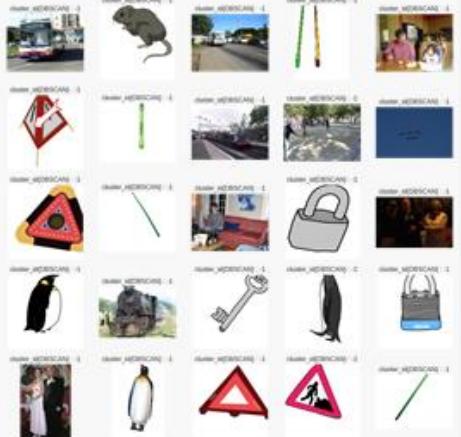
NIA 한국지능정보사회진흥원

DBSCAN

Solution 2 - Clustering

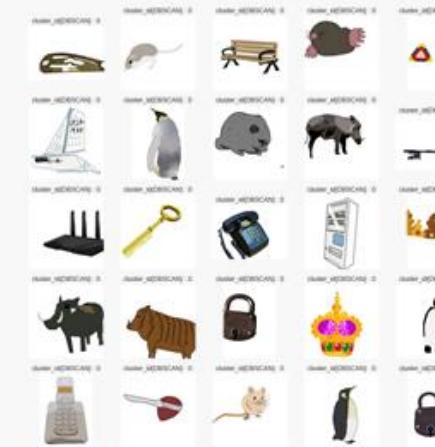
DBSCAN 결과

Example of id -1 dataset clustered by DBSCAN



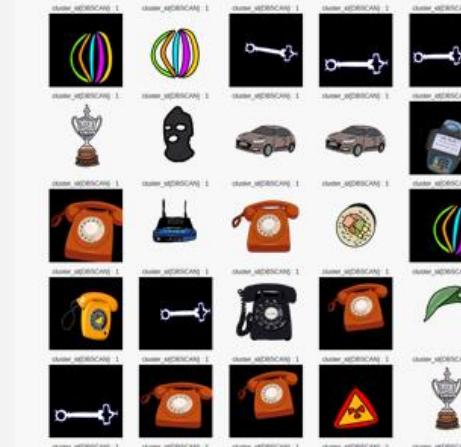
Cluster -1

Example of id 0 dataset clustered by DBSCAN



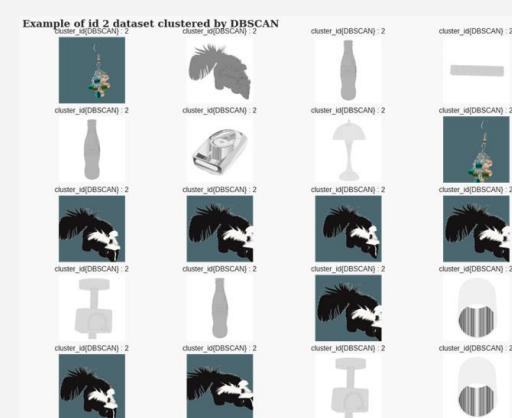
Cluster 0

Example of id 1 dataset clustered by DBSCAN



Cluster 1

Example of id 2 dataset clustered by DBSCAN



Cluster 2

실사 이미지들은
Cluster 되지 못함



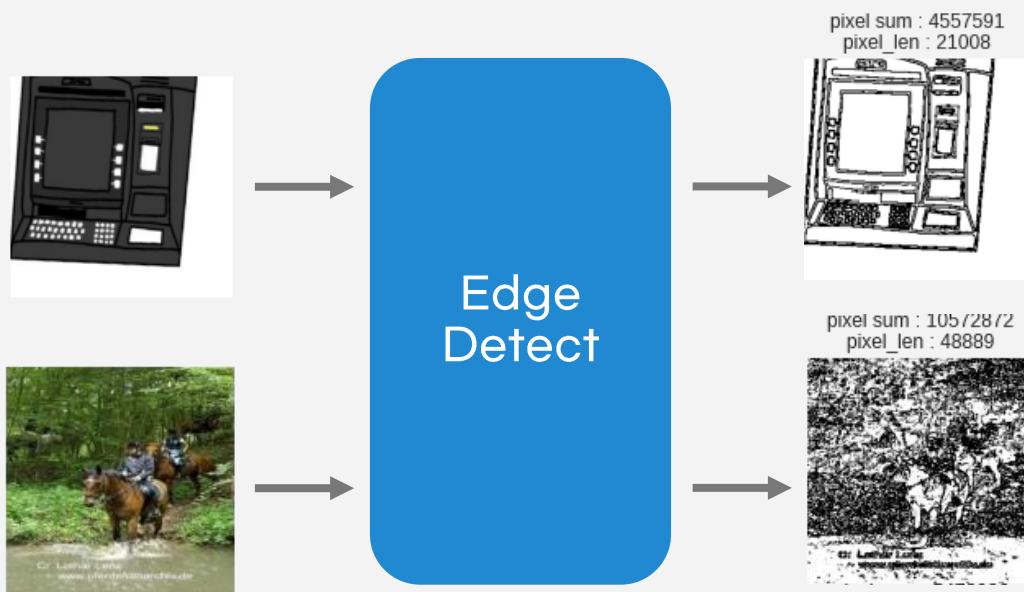
과학기술정보통신부

NIA 한국지능정보사회진흥원

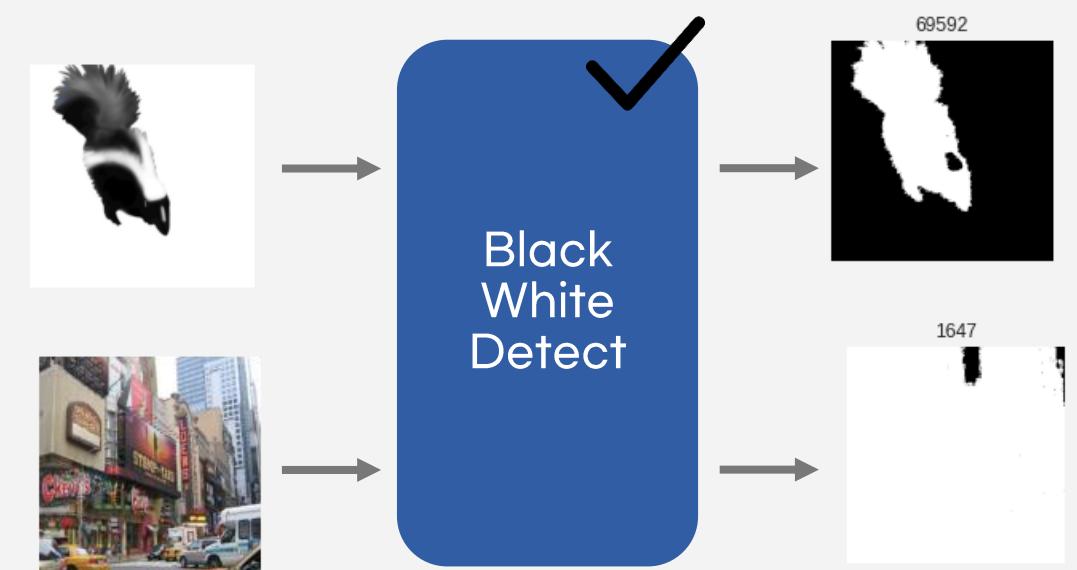
필터링

Solution 2 - Clustering + Filtering

Edge Detection Filtering



Black White Filtering

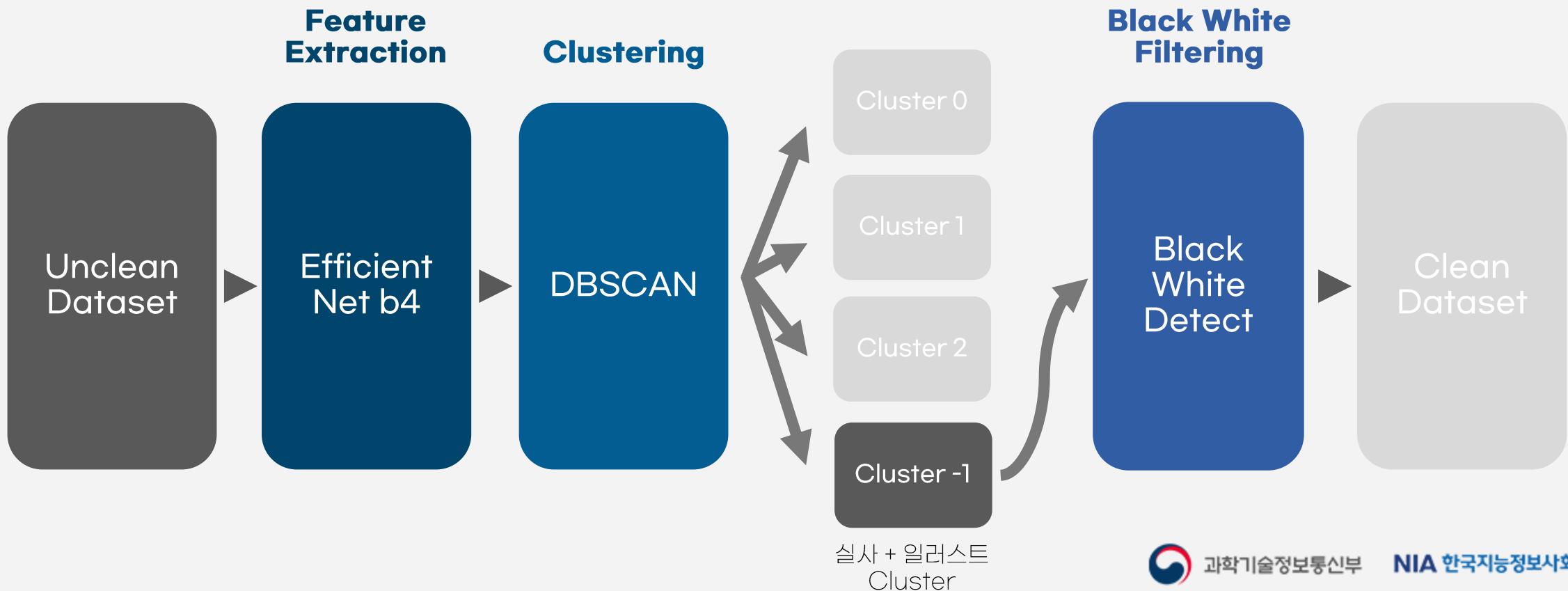


과학기술정보통신부

NIA 한국지능정보사회진흥원

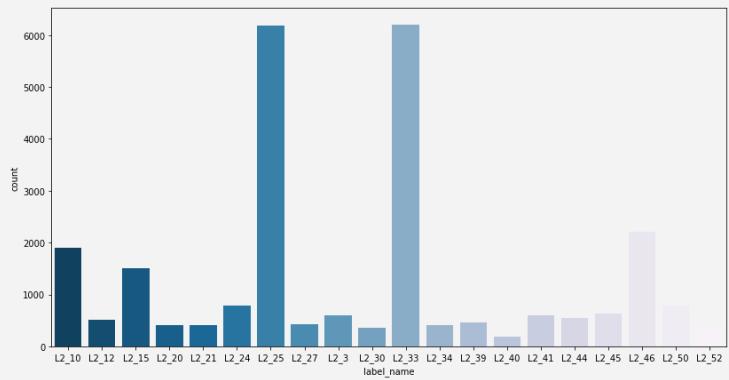
최종 실사 이미지 제거

Solution 2 - Clustering + Filtering



데이터 불균형 문제 해결

클래스별 데이터 분포



원본
데이터셋

25,300
→
23,000

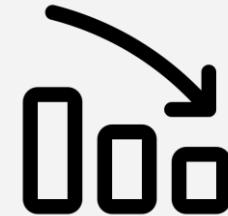
실사 이미지
제거후

Imbalanced Data 해결 아이디어



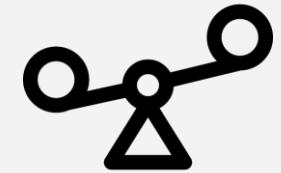
Oversampling

이미지 증강



Undersampling

랜덤 샘플링



Class Weight

가중치 적용



과학기술정보통신부

NIA 한국지능정보사회진흥원

Oversampling

이미지 증강

회전, 두집기, 대비/밝기변화, 이미지크기변화, 크롭

```
def albu_transforms():
    return A.Compose([
        A.Rotate(),
        A.OneOf([
            A.HorizontalFlip(p=1.0),
            A.VerticalFlip(p=1.0),
            ], p=0.5),
        A.OneOf([
            A.RandomContrast(limit=0.5, p=1.0),
            A.RandomGamma(gamma_limit=(200, 250), p=1.0),
            A.RandomBrightnessContrast(brightness_limit=0.4, contrast_limit=0.4, p=1.0),
            ], p=1.0),
        A.OneOf([
            A.Resize(256, 256),
            A.Compose([
                A.Resize(320, 320),
                A.CenterCrop(256, 256),
            ])
        ], p=1.0)
    ])
```

클래스별 증강 배수 설정

데이터가 매우 적은 클래스
L2_12

× 3

데이터가 적은 클래스
L2_3, L2_40, L2_41

× 2



과학기술정보통신부

NIA 한국지능정보사회진흥원

Undersampling

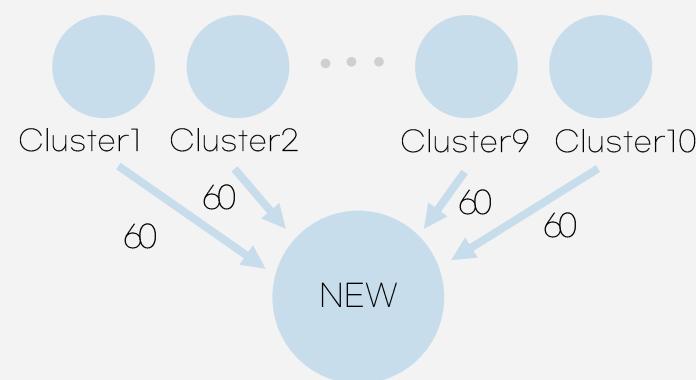
데이터가 평균 이상인 클래스

L2_15, L2_46, L2_25, L2_10, L2_33



Solution1-Clustering

K=10군집에서 군집별 60개씩 데이터샘플링



Solution2-Random sampling

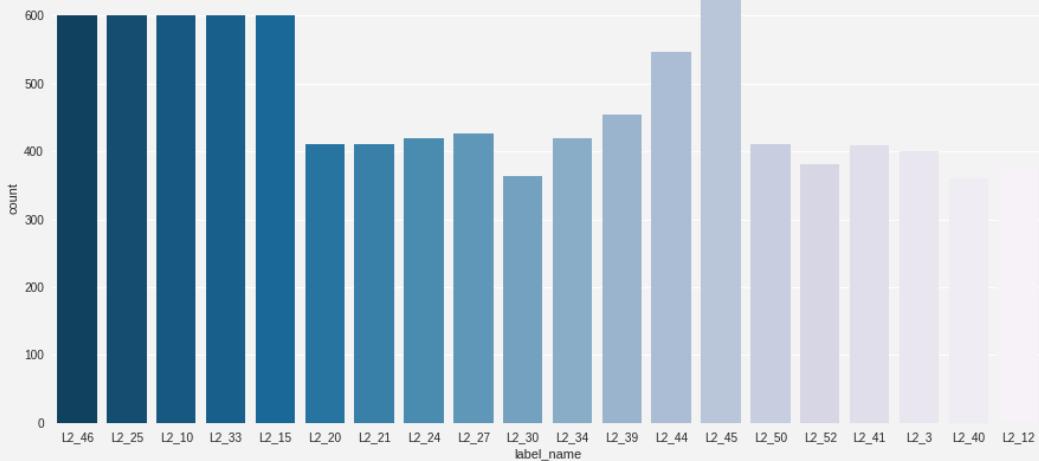
클래스별 랜덤 600개씩 데이터샘플링



과학기술정보통신부

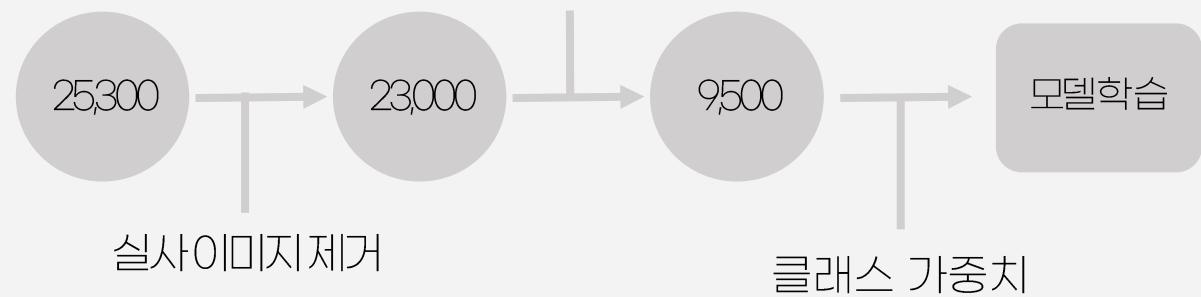
NIA 한국지능정보사회진흥원

클래스 가중치



불균형 문제 해결 후 데이터 분포 변화

Oversampling
&
Undersampling



모델 학습 시 클래스 가중치 적용

```

class_weights = compute_class_weight(class_weight = "balanced",
                                      classes = np.unique(df['label_name']),
                                      y = df['label_name'])
class_weights = dict(enumerate(class_weights))

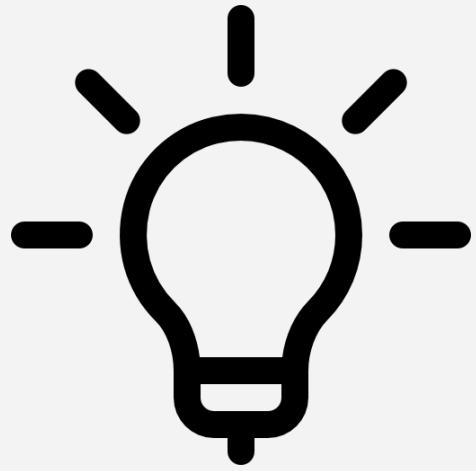
history = model5.fit_generator(train_dataset, validation_data = validation_dataset,
                                epochs=n_epoch, class_weight = class_weights, callbacks = [e
    
```



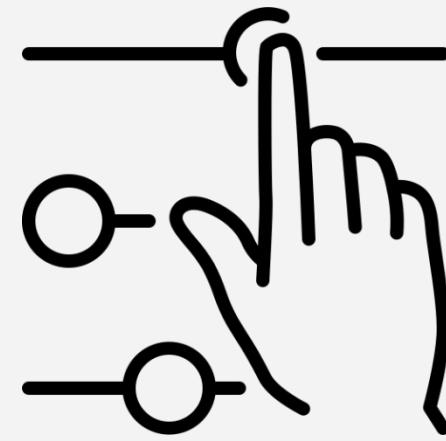
과학기술정보통신부

NIA 한국지능정보사회진흥원

모델 설계

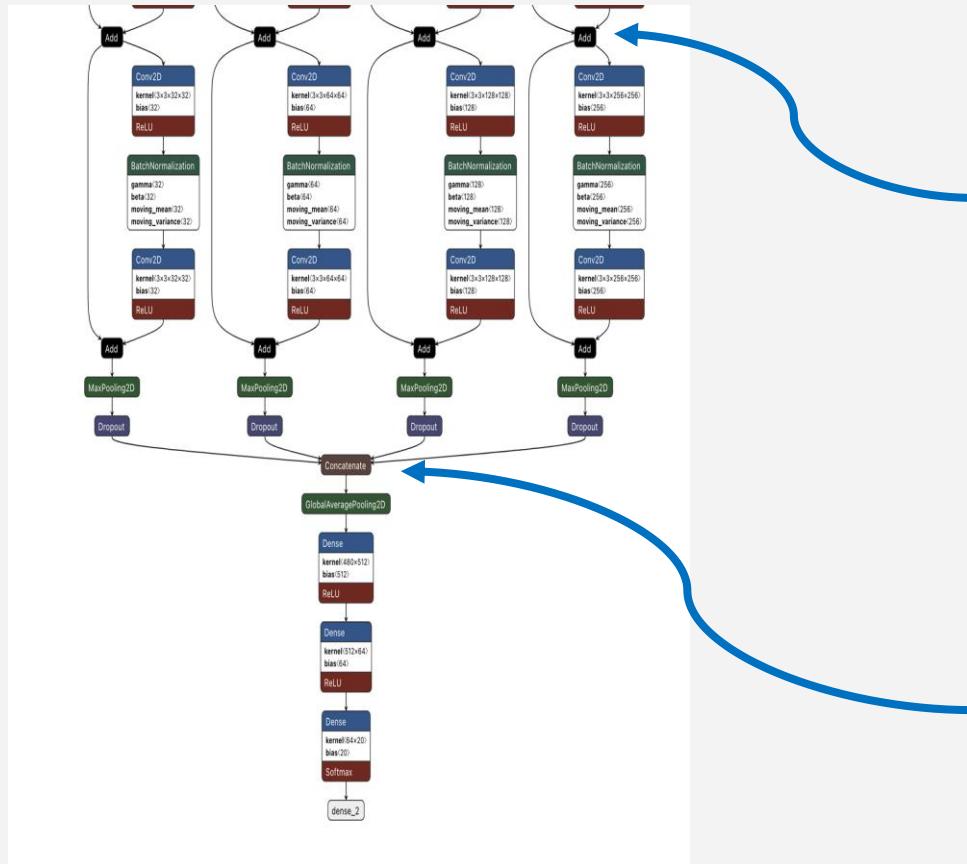


①
ResNet,
InceptionNet
아이디어 모델

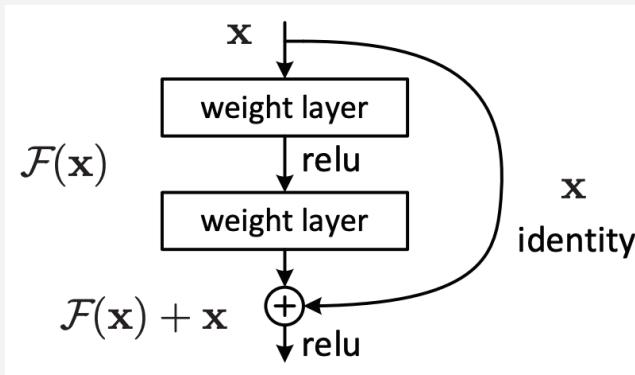


②
CNN
커스텀 모델

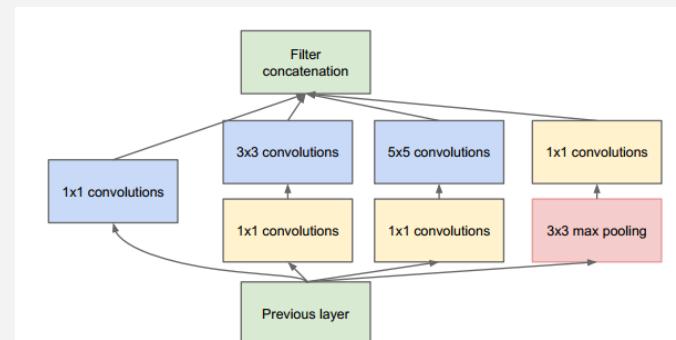
ResNet, InceptionNet 아이디어



Residual Block



Inception Block



잔여효과로 인해
학습이 더 쉬워지고
기울기 소실 문제 완화

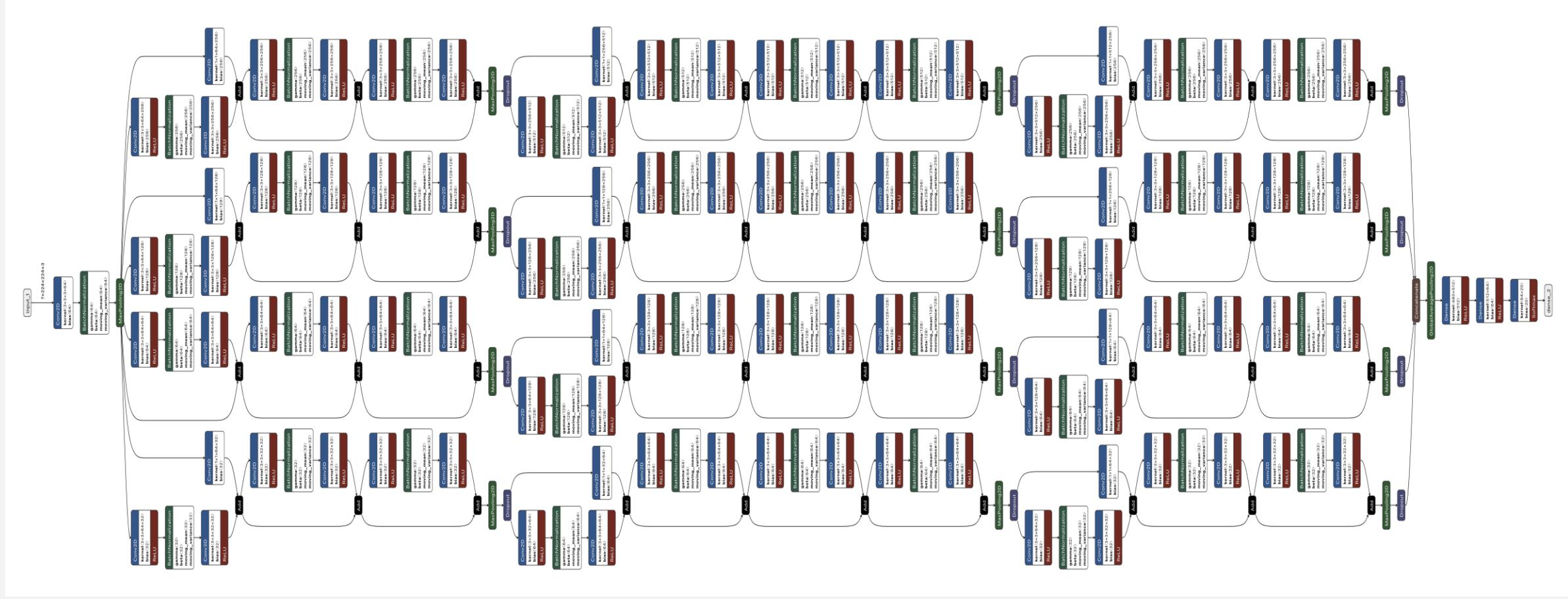
다양한 채널에서
추출된 Feature를
결합할 수 있음



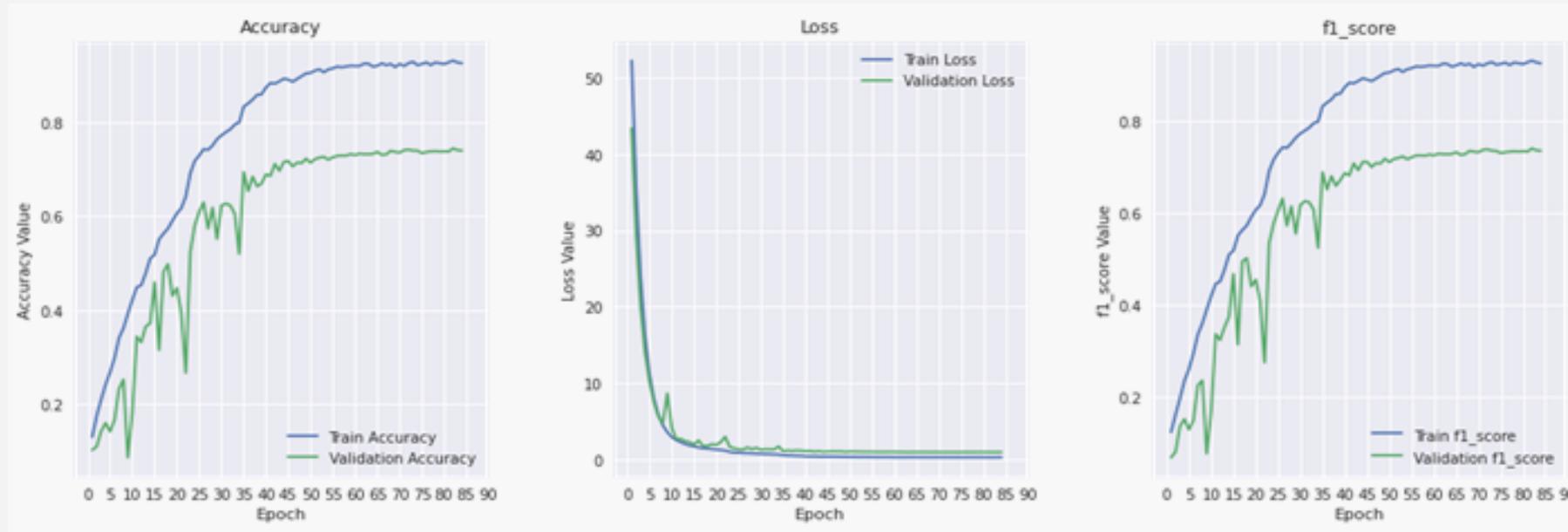
과학기술정보통신부

NIA 한국지능정보사회진흥원

ResNet, InceptionNet 아이디어 모델



ResNet, InceptionNet 아이디어 모델 성능평가



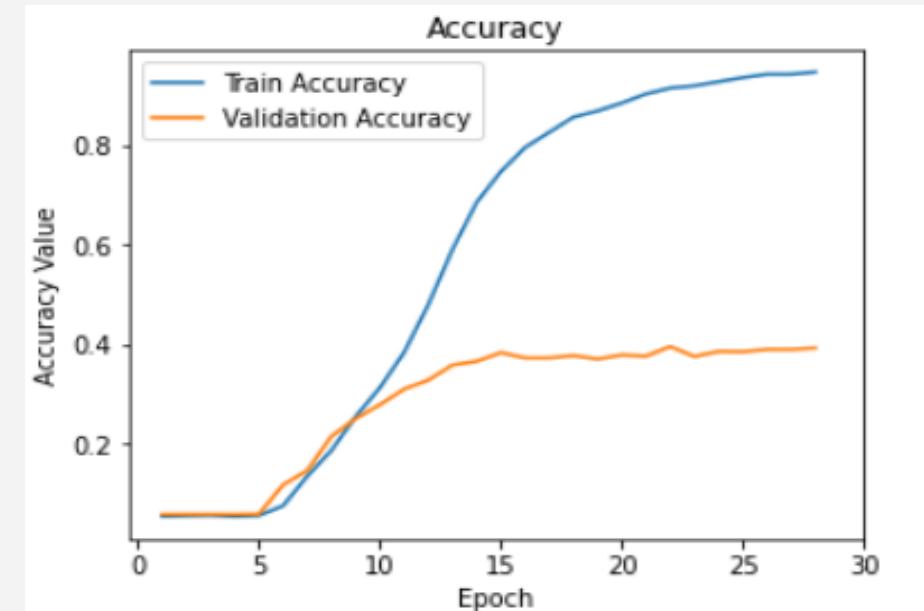
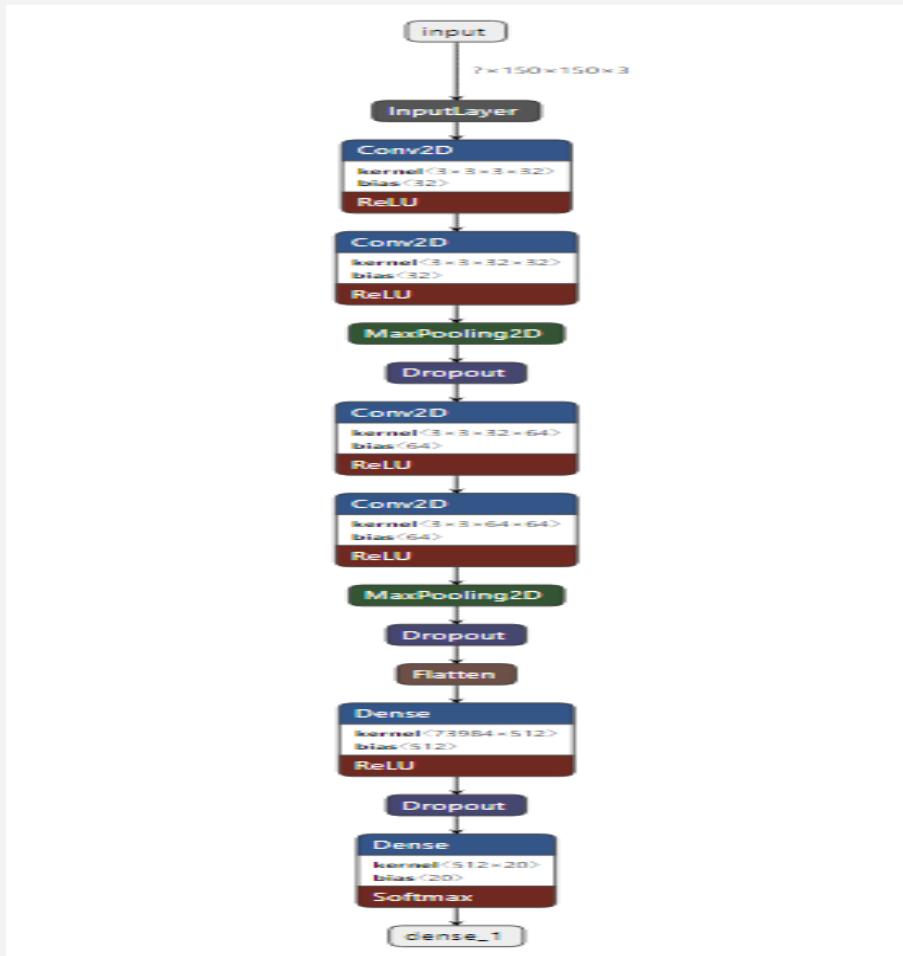
F1 score = 0.7337



과학기술정보통신부

NIA 한국지능정보사회진흥원

CNN 커스텀 모델 - Basic



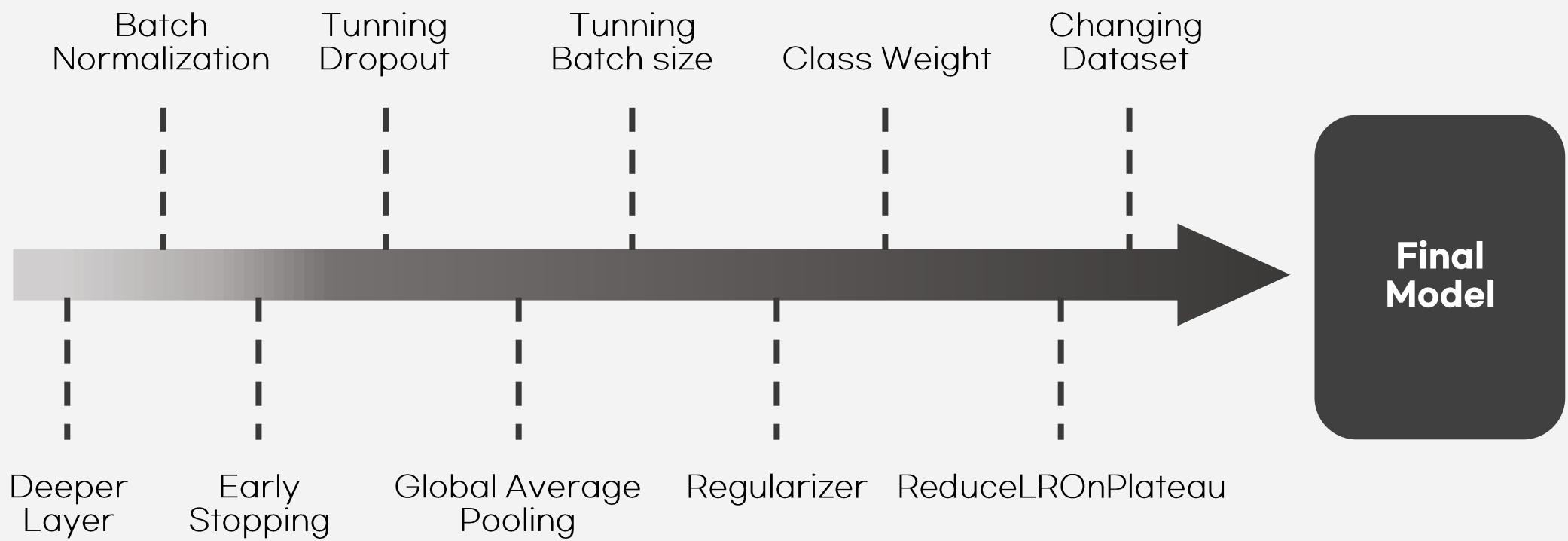
CNN 커스텀 모델 실험 전
기준 모델 생성



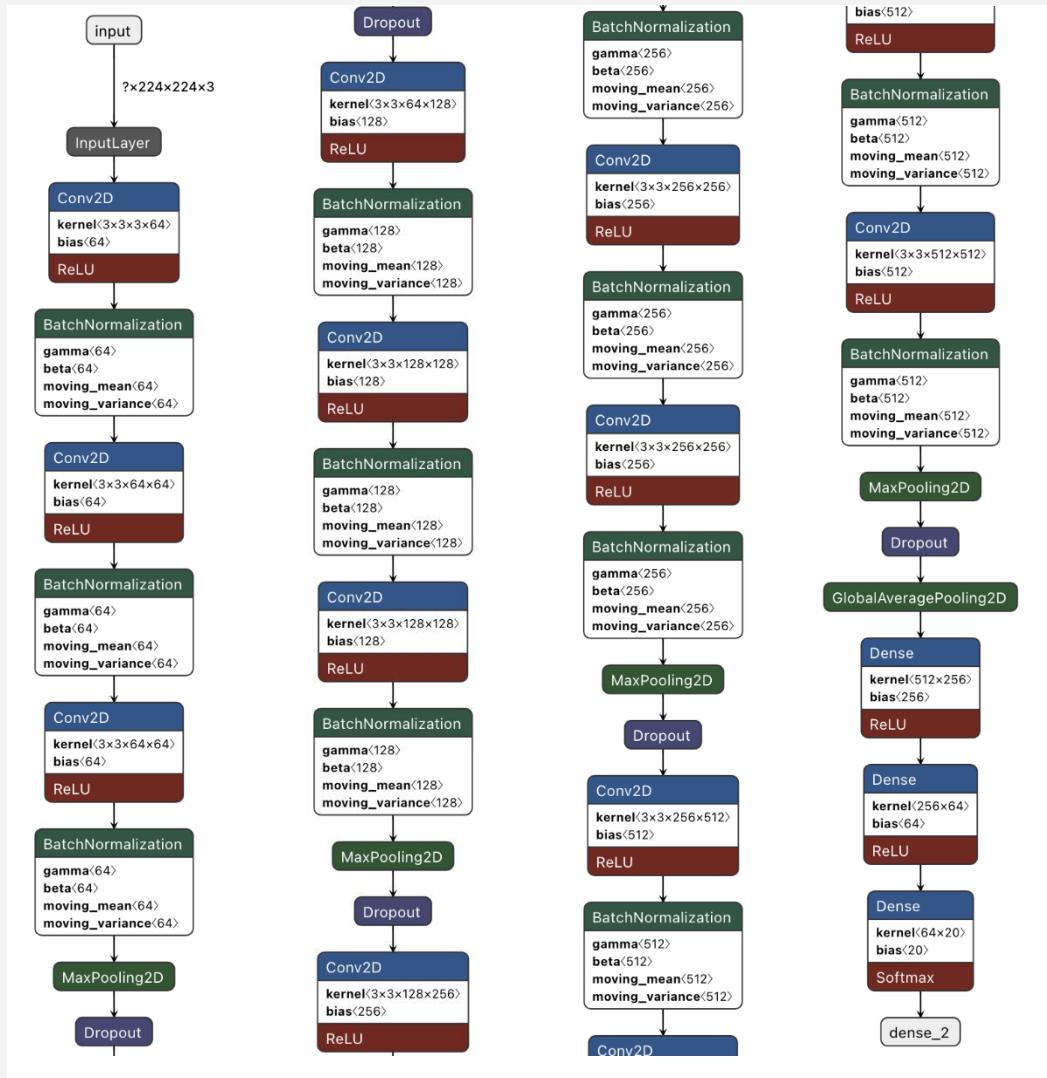
과학기술정보통신부

NIA 한국지능정보사회진흥원

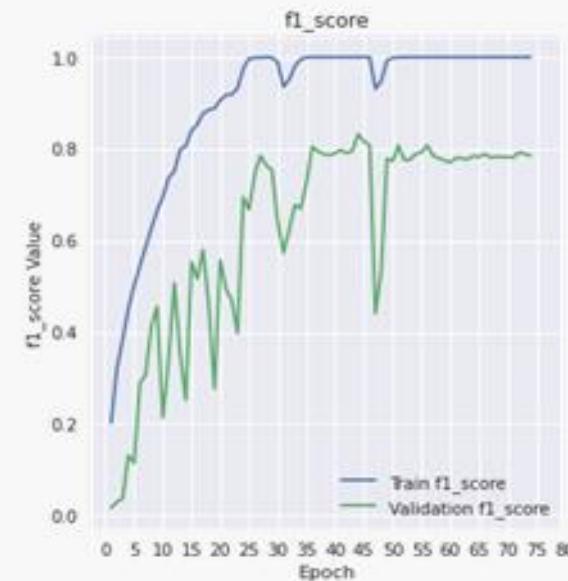
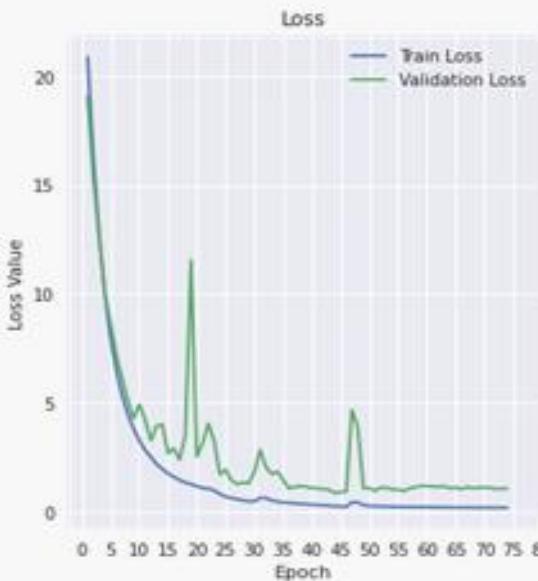
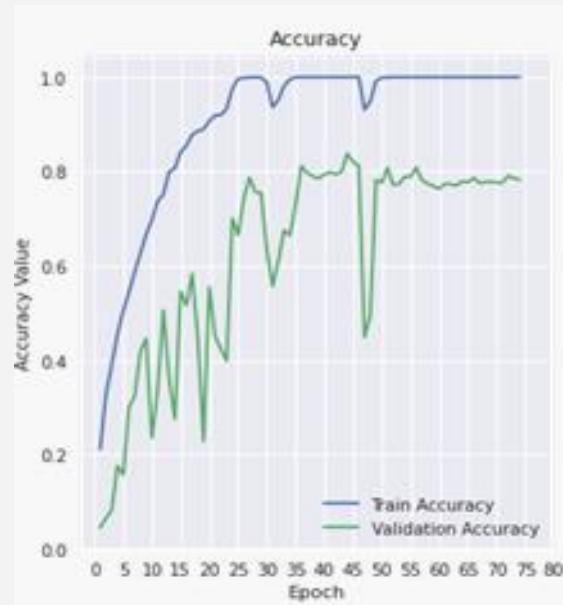
CNN 커스텀 모델



CNN 커스텀 모델 - Final



CNN 커스텀 모델 성능평가



F1 score = 0.8317



과학기술정보통신부

NIA 한국지능정보사회진흥원



감사합니다

2022 DATA CREATOR CAMP



과학기술정보통신부

NIA 한국지능정보사회진흥원