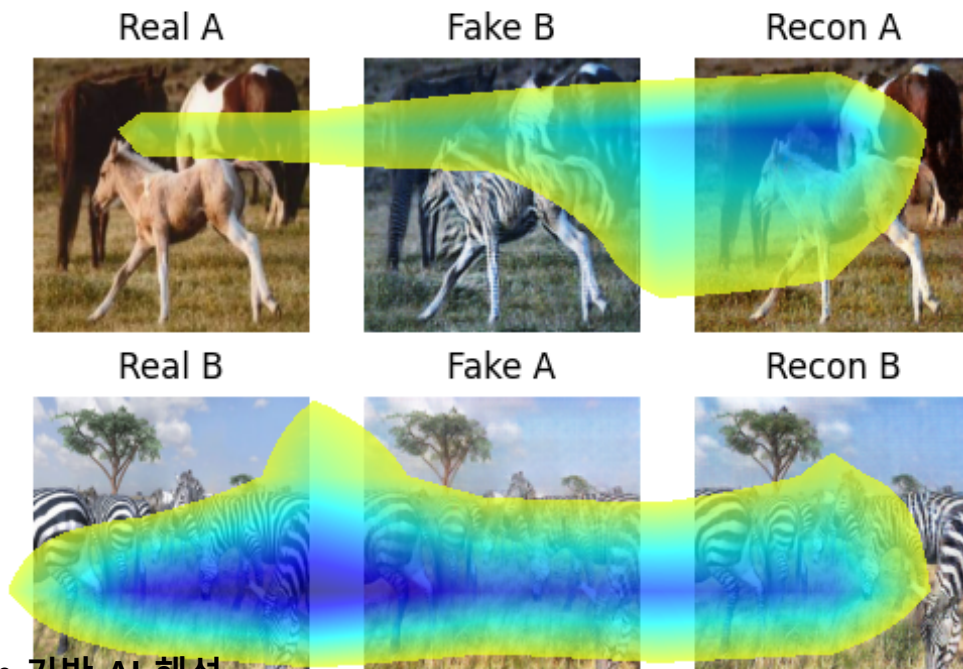


딥페이크 히트맵 분석 보고서

1. 분석 개요

- 모델명: MobileNetV3-Small
- 모델 유형: 한국인 전용 모델
- 분석 일시: 2025-11-10 00:19:50
- 예측 결과: Fake (95.08%)
- 딥페이크 확률: 29.30%

2. Grad-CAM 시각화



3. LangChain 기반 AI 해석

Grad-CAM(Gradient-weighted Class Activation Mapping) 히트맵은 딥러닝 모델의 판단을 시각적으로 해석하는데 매우 유용한 도구입니다. 주어진 정보에 기반하여 예측 결과의 신뢰도 및 한계점, 그리고 중요한 시각적 요소들을 다음과 같이 분석할 수 있습니다.

Grad-CAM 해석

1. **붉은색 영역**:

- 붉은색으로 표시된 영역은 모델이 "Fake" 판단을 내리는데 있어 주요 근거로 삼았던 부분입니다. 이는 모델이 주의 깊게 살펴본 부분이며, 합성된 이미지의 특성을 파악하는 데 도움이 됩니다. 일반적으로 딥페이크 이미지에서는 합성 흔적, 피부 질감, 조명 왜곡이 주요 요소로 작용할 수 있습니다.

2. ****합성 흔적****:

- 딥페이크 기술은 종종 이미지의 가장자리에서 어색한 경계선이나 불연속성을 나타냅니다. 예를 들어, 얼굴과 배경이 자연스럽게 통합되지 않을 경우, 모델은 이러한 비대칭성을 감지할 수 있습니다. 붉은색 영역에서 나타나는 갑작스러운 변화는 합성을 나타낼 수 있으며, 이는 모델이 해당 부분에서 눈에 띄는 이상을 발견했음을 시사합니다.

3. ****피부 질감****:

- 자연스러운 피부 텍스처는 사람의 얼굴에서 특유의 미세한 디테일(예: 주름, 여드름, 모공)이 포함되는 반면, 합성된 이미지는 이러한 디테일이 결여되거나 비정상적으로 부자연스럽게 나타날 수 있습니다. 붉은색 히트맵이 이러한 피부 질감의 부조화에 주목했다면, 이는 딥페이크라는 판단을 강화하는 요소가 될 수 있습니다.

4. ****조명 왜곡****:

- 실시간 비디오나 이미지에서의 조명은 얼굴의 윤곽과 느낌을 크게 좌우합니다. 딥페이크의 경우, 조명이 얼굴에 고르게 적용되지 않거나 주변의 조명과 맞지 않아 어색한 그림자나 빛 반사가 발생하는 경우가 많습니다. 모델이 이러한 조명 왜곡을 감지했을 가능성이 있습니다.

신뢰도와 한계점

1. ****신뢰도****:

- 예측 결과가 "Fake"라는 높은 확률(95.08%)로 나타난 것은 모델의 학습 데이터가 잘 구성되어 있어 해당 이미지를 신뢰성 있게 분석했음을 나타냅니다. 특히 Grad-CAM 히트맵을 통해 중요한 시각적 요소들을 확인할 수 있기 때문에 예측의 해석 가능성이 높아집니다.

2. ****한계점****:

- 그러나 Grad-CAM의 해석은 주관적일 수 있으며, 붉은색 영역이 실제로 모델이 고려한 모든 요소를 반영한다고 보장할 수는 없습니다. 즉, 딥페이크와 같은 복잡한 이미지가 항상 명확한 신호를 포함하지 않을 수 있으며, 이는 오탐률(False Positive)을 증가시킬 수 있습니다. 또한, 노이즈 많은 환경이나 저해상도 이미지에서는 모델의 성능이 저하될 수 있습니다.

3. ****심층 결과****:

- 딥페이크 탐지의 심층 결과는 다양한 세부 사항을 포함할 수 있습니다. 예를 들어, 모델이 사용한 특정 특징 추출 방식이나 전처리 단계가 결과에 미친 영향을 이해하는 것이 중요합니다. 또한, 데이터셋의 편향(ex. 특정 인종, 성별, 조명 조건)도 영향을 줄 수 있습니다. 따라서 다양한 환경에서 모델을 정확하게 평가하고, 다양한 유형의 합성 이미지를 테스트하는 것이 필요합니다.

결론적으로, Grad-CAM은 모델의 판단을 설명하는 유용한 도구이지만, 해석 시 주의가 필요하며, 지속적인 연구와 검증이 필수적입니다.

40. 결론 및 권장 조치

본 분석은 Grad-CAM 시각 주목도를 중심으로 진행되었습니다.

AI의 결과는 참고용으로 사용해야 하며, 법적 판단이나 공식 증거로 사용되지 않습니다.

결과의 신뢰도를 높이기 위해 다양한 이미지 소스로 교차 검증을 권장합니다.