

딥페이크 히트맵 분석 보고서

1. 분석 개요

- 모델명: MobileNetV3-Small
- 모델 유형: 한국인 전용 모델
- 분석 일시: 2025-11-05 14:22:32
- 예측 결과: 이 이미지는 Real로 분류되었으며, 예측신뢰도는 98.11%입니다.
- 딥페이크 확률: 23.03%

2. Grad-CAM 시각화



3. LangChain 기반 AI 해석

Grad-CAM(Gradient-weighted Class Activation Mapping) 히트맵을 기반으로 분석한 내용을 다음과 같이 정리할 수 있습니다.

Grad-CAM 히트맵 해석

1. **예측 결과 해석**

- 모델의 예측 결과에 따르면, 해당 이미지는 **Real**로 분류되었으며, **신뢰도 98.11%**라는 높은 수치를 보이고 있습니다. 이는 모델이 해당 이미지가 진짜일 가능성이 매우 높다고 판단했다는 것을 의미합니다.

- 하지만, **딥페이크 확률 23.03%**는 어느 정도의 불확실성이 존재함을 나타내며, 이론적으로는 해당

이미지가 합성일 가능성도 있다는 것을 시사합니다.

2. **붉은색 영역의 의미**

- **붉은색 영역**: 히트맵에서 붉은색으로 강조된 부분은 모델이 딥페이크 여부를 판단하는 데 있어 중요한 시각적 요소로 인식한 부분입니다. 이는 보통 이미지의 피부 질감, 조명 변화, 그리고 합성의 흔적이 강하게 나타나는 부위일 가능성이 높습니다.

- **피부 질감**: 딥페이크 이미지에서는 종종 피부 질감이 비정상적으로 매끈하거나 부자연스러운 경우가 있습니다. 붉은색 영역이 얼굴의 특정 부분에 집중되어 있다면, 그 해당 부위의 질감이 비정상적이거나 인위적으로 조정되었을 가능성이 있습니다.

- **조명 왜곡**: 또한 조명의 일관성이 부족한 경우도 딥페이크의 흔적인데, 붉은색으로 강조된 부분이 일관된 방향에서 오는 조명을 따르지 않거나 그림자 처리에 문제가 있을 경우 의심스러운 요소로 기능할 수 있습니다.

신뢰도와 한계점

- **신뢰도**: 98.11%라는 높은 신뢰도는 모델의 전반적인 성능을 나타내지만, 이 수치가 '완벽함'을 보장하지는 않습니다. 딥페이크 탐지 모델은 다양한 변수를 고려해야 하며, 특정 이미지나 데이터셋에 따라 성능이 달라질 수 있습니다.

- **한계점**:

- **과신 위험**: 모델이 높은 신뢰도를 보일 경우 진짜 이미지가 아닌 경우에도 잘못된 결정을 내릴 위험이 있습니다. 예를 들어, 실제 사람의 사진에서 특이한 조명이나 각도에 의해 의심스러운 특징이 나타날 수 있습니다.

- **시각적 편 포인트 부족**: Grad-CAM 기반 분석은 시각적으로 중요한 구역을 시각화하는 데 유용하지만, 이미지의 전반적인 품질 변화나 복잡한 합성 패턴은 설명하기 어려울 수 있습니다.

- **환경적 요인**: 촬영된 환경의 차이, 조명 조건, 배경의 복잡성이 모델의 성능에 영향을 줄 수 있습니다.

추가 심층 결과

- **세부적인 특징 검토**: 붉은색 영역 외에도 명확히 드러나지 않는 미세한 결함이나 불일치가 있을 수 있으며, 이러한 부분은 인간 전문가의 눈으로 더욱 체계적이고 상세한 검토가 필요합니다.

- **다양한 데이터셋의 필요성**: 모델의 신뢰도를 높이기 위해서는 다양한 데이터셋을 활용하여 훈련시켜야 하며, 실제 환경에서의 변화를 반영한 데이터를 포함하는 것이 중요합니다.

결론적으로, Grad-CAM 히트맵은 딥페이크 탐지에 중요한 도구가 될 수 있으나, 모델의 결정에 있어 인간의 추가적인 분석과 판단이 필요합니다.

4. 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.