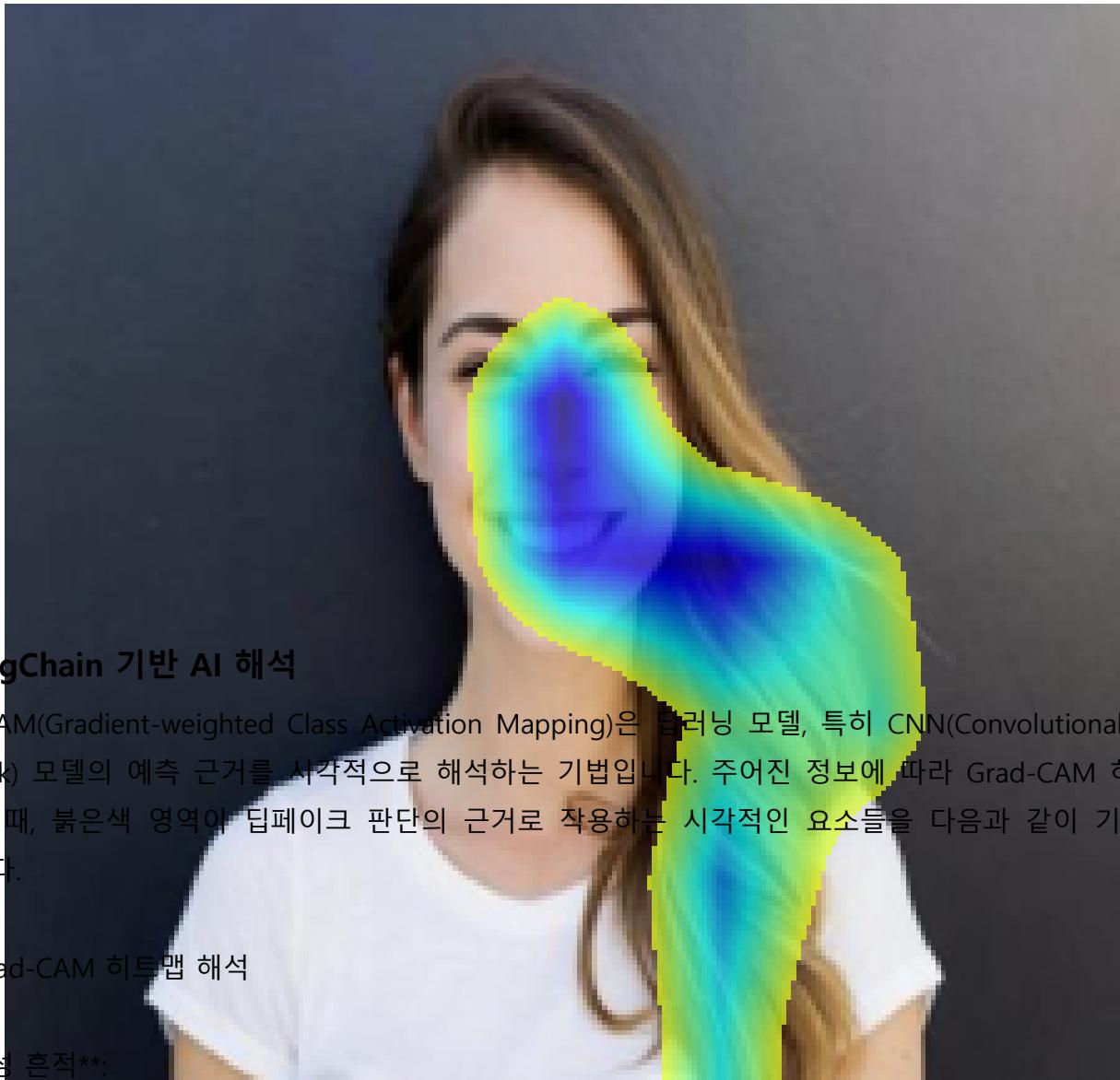


딥페이크 히트맵 분석 보고서

1. 분석 개요

- 모델명: MobileNetV3-Small
- 모델 유형: 외국인 전용 모델
- 분석 일시: 2025-11-06 15:20:44
- 예측 결과: Real (7.69%)
- 딥페이크 확률: 14.70%

2. Grad-CAM 시각화



3. LangChain 기반 AI 해석

Grad-CAM(Gradient-weighted Class Activation Mapping)은 딥러닝 모델, 특히 CNN(Convolutional Neural Network) 모델의 예측 근거를 시각적으로 해석하는 기법입니다. 주어진 정보에 따라 Grad-CAM 히트맵을 분석할 때, 붉은색 영역이 딥페이크 판단의 근거로 작용하는 시각적인 요소들을 다음과 같이 기술할 수 있습니다.

Grad-CAM 히트맵 해석

1. **합성 흔적**:

- 붉은색으로 표시된 영역에서 합성 흔적이 눈에 띕니다. 이 영역은 자연스러운 피부 질감이나 주름이 없는 부위로, 인조적인 평면감이 돌보일 수 있습니다.

- 합성된 이미지에서 자주 나타나는 경계선이 부자연스럽거나 날카로운 경우, 이 부분은 딥페이크 모델이 인지하는 주요 특징으로 볼 수 있습니다.

2. **피부 질감**:

- 자연스러운 피부 질감과 딥페이크 이미지 간의 눈에 띄는 차이가 존재할 경우, 이러한 영역은 붉은색으로 강조됩니다.

- 일반적으로 진짜 이미지에서는 세밀한 주름이나 결이 보이고, 딥페이크에서는 내용이 매끄럽고 반복적인 패턴일 수 있습니다. 이로 인해 붉은 지역이 더 두드러지게 나타날 수 있습니다.

3. **조명 왜곡**:

- 모델이 인식한 붉은색 영역 중 조명과 그림자가 불균형한 경우도 중요한 탐지 요소로 작용합니다.

- 자연스러운 이미지에서는 조명 분포가 얼굴의 윤곽에 맞춰 다이내믹하게 변하는 반면, 딥페이크에서는 인공적인 조명 처리가 이루어질 수 있고, 그 결과 불균형한 색도우가 나타나는 경우가 많습니다.

신뢰도 및 한계점

- **신뢰도**: Grad-CAM이 제공하는 시각적 설명은 모델이 특정 픽셀이나 영역에 얼마나 의존하는지 보여 줍니다. 이는 모델의 결정 과정에 대한 심층적인 이해를 가능하게 하여, 예측이 왜 발생했는지를 설명하는 데 유용합니다. 특히, 흉터처럼 보이는 합성 흔적이나 불일치한 피부 질감은 고전적인 인간의 평론가들이 딥페이크를 식별하는 데 유용한 팁이 될 수 있습니다.

- **한계점**: 그러나 Grad-CAM의 해석은 주관적일 수 있으며, 모든 경우에 정확한 해석을 보장하지 않습니다. 예를 들어, 모델이 잘못 추론한 경우 낮은 신뢰도의 히트맵이 생성될 수 있습니다. 또한, 특정 데이터 세트에서 훈련된 모델의 경우, 비슷한 유형의 이미지에 대해 과적합(overfitting) 문제로 잘못된 판단을 할 수 있습니다.

추가 심층 결과

1. **기타 특징**:

- 피부 톤의 불일치: 진짜 얼굴과 합성된 얼굴 간의 톤 차이를 감지할 수 있습니다. 이 경우, 대부분의 모델이 감지하기 어려운 미세한 차이들이 존재할 수 있습니다.

- 눈의 움직임: 눈의 각도나 iris의 생동감 차이 또한 모델이 주목할 수 있는 부분입니다.

2. **결론적으로**:

- Grad-CAM 히트맵은 모델이 어디에서 딥페이크의 허점을 찾아내는지에 대한 유용한 통찰을 제공하지만, 심층적인 탐지를 위해서는 여러 기법의 결합이 필요할 수 있습니다. 예를 들어, GAN(Generative Adversarial Network)에서 생성된 특징 분석과 함께 사용하면 더욱 강화된 탐지 성능을 기대할 수 있습니다.

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.