

딥페이크 히트맵 분석 보고서

1. 분석 개요

- 모델명: MobileNetV3-Small
- 모델 유형: 외국인 전용 모델
- 분석 일시: 2025-11-07 09:20:42
- 예측 결과: Fake (52.49%)
- 딥페이크 확률: 18.80%

2. Grad-CAM 시각화



3. LangChain 기반 AI 해석

Grad-CAM(Gradient-weighted Class Activation Mapping)은 신경망 모델이 특정 예측을 내리기 위해 어떤 부분에 주목했는지를 시각적으로 나타내는 기법으로, 모델의 해석 가능성을 높이는 데 유용합니다. 주어진 정보를 바탕으로, 예측 결과가 "Fake"로 나왔고 Grad-CAM 히트맵에서 붉은색 영역이 강조된 것에 대해 분석해보겠습니다.

시각적 근거 분석

1. **합성 흔적 (Artifacts):

- 붉은색 영역에서 합성 흔적을 발견한 경우, 이는 일반적으로 이미지의 경계, 특히 얼굴 윤곽이나 턱선, 머리카락의 경계에서 나타나는 부자연스러운 픽셀들입니다. 이런 지역은 종종 실제 이미지는 아닌

인공적으로 합성된 이미지에서 나타나는 영향이 있습니다. 모델이 이러한 영역에 높은 관심을 기울였다면, 이는 합성이 이루어진 증거로 판단할 수 있습니다.

2. **피부 질감**:

- 비정상적인 피부 질감은 또 다른 중요한 요소입니다. 딥페이크 기술이 피부의 질감을 조작할 때 종종 자연스러운 피부의 결점이나 주름 없이 너무 고르게 만들어질 수 있습니다. 모델이 피부 질감에 주목했다면, 너무 매끄럽거나 비정상적으로 결이 없는 부분에서 경고 신호를 인식했을 가능성이 큽니다. 이러한 차이로 인해 모델은 딥페이크라고 판단했을 수 있습니다.

3. **조명 왜곡**:

- 조명은 자연스럽고 일관된 패턴을 따라야 하는데, 인공적으로 생성된 이미지에는 조명이 비정상적으로 배치되거나, 그림자가 불균형하게 나타날 수 있습니다. 붉은색 영역이 특정 조명 효과를 보인다면, 이는 합성이 이루어진 부분에서 조명 처리가 부자연스럽다는 것을 나타내며, 따라서 모델이 이를 딥페이크의 신호로 인식했을 수 있습니다.

신뢰도 및 한계점

1. **신뢰도**:

- Grad-CAM은 시각적으로 직관적으로 모델의 판단 과정을 이해하는 데 큰 도움을 줍니다. 특히, 특정 영역이 모델의 결정에 미친 영향을 명확하게 보여주어 연구자나 저널리스트가 깊이 있는 분석을 할 수 있게 합니다. 모델이 특정 이상 징후를 발견하면서 딥페이크가 적어도 52.49% 확률로 의심스럽다고 판단한 것은 주목할 만한 결과입니다.

2. **한계점**:

- 그러나 Grad-CAM의 해석은 비단 이미지의 시각적 요소만으로 모든 판단을 설명할 수는 없습니다. 모델이 붉은 영역에 배정한 중요도가 모델의 학습 데이터에 의해 좌우될 수 있으며, 인위적으로 합성된 이미지와 진짜 이미지의 경계를 명확히 할 수 있는 유일한 기준이 되지 않습니다. 또한, 일반적으로 과적합된 모델이 덜 일반화된 특성을 가져서 비정상적인 결과를 초래할 수 있습니다.

심층 결과

추가적으로 모형이 대용량 데이터로 학습된 경우, 감지의 정확도는 상대적으로 높아질 수 있으나, 그러한 데이터가 실제로 합성된 이미지의 다양성을 충분히 포괄하지 않는다면, 오탐지율과 미탐지율이 동시에 증가할 수 있습니다. 즉, 모델이 훈련한 데이터 분포와 실제 환경에서 나타나는 분포의 불일치로 인해 신뢰도가 낮아질 수 있습니다.

결론적으로, 이 히트맵은 특히 특정한 시각적 요소의 결함이 딥페이크 판단에 중요한 역할을 하고 있다는 것을 나타내지만, 이를 절대적인 신뢰의 기준으로 삼아서는 안 됩니다.

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.