

生成モデルを用いたニューラルネットワーク 検証における検証領域近似手法の改良

富士通株式会社

横山晴樹

近似を用いたNN検証の問題点

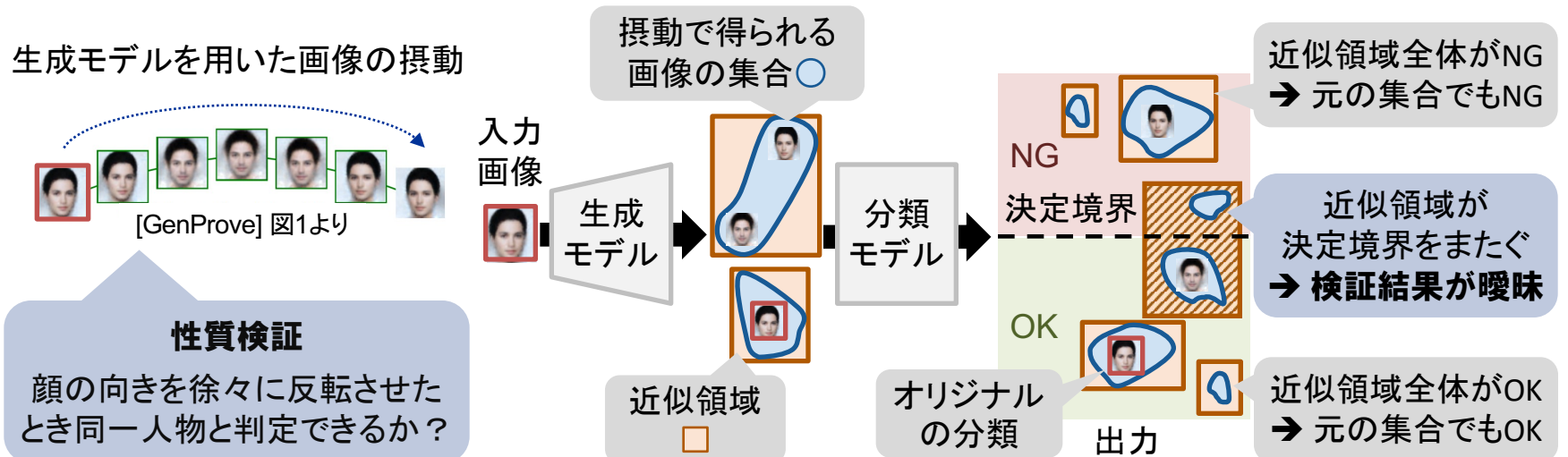
- 「入力に摂動を与えてもNNが誤認識しないか」等の性質を検証する技術では、スケーラビリティのため検証する領域を近似する
- 性質検証において、近似領域が決定境界を跨いでいると、真偽が曖昧となる

手法の改良

少数のサンプルに対するNNの出力を用いて決定境界からの距離を推定し、近似領域の大きさを制御する

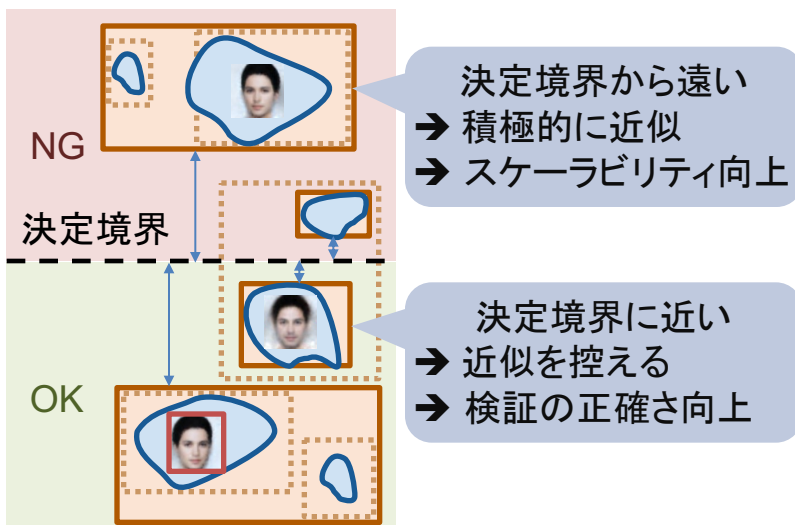
→ 決定境界を跨ぐ近似領域の削減

近似を用いたNN検証



改良点

少数サンプルに対し決定境界からの距離を測定
→ 距離に応じた近似戦略



評価

