

유전체 친족 관계 분류 연구 및 모델 최적화 최종 결론 보고서

1. 분석 성과 및 주요 지표 (Performance Metrics)

- 최고 정확도: **cM_1** 데이터셋 기준 **96.50%**의 정확도를 기록하며 초기 벤치마크 목표를 완수하였습니다.
- 데이터 스펙트럼의 효용성: 분석 결과, 세그먼트 임계값을 1cM으로 낮추어 짧은 유전적 세그먼트 정보를 최대한 보존하는 것이 분류 성능 향상의 핵심 동인임을 확인하였습니다. (cM_10 대비 정확도 약 4%p 향상)
- 모델별 랭킹 (Global Leaderboard):
 - LightGBM (96.50%)**: 분석 효율 및 정확도 분문 최적 솔루션
 - CatBoost (95.53%)**: 과적합 방지 및 모델 안정성 우수
 - XGBoost (95.21%)**: 빠른 학습 속도 및 안정적인 성능 기록
 - RandomForest (95.02%)**: 전통적 분류 기법 기반의 견고한 기준점 제공

2. 모델 신뢰성 및 검증 (Validation & Anti-Overfitting)

높은 정확도가 특정 데이터에 매몰된 결과(과적합)가 아님을 입증하기 위해 다음과 같은 검증 프로토콜을 수행하였습니다.

- 5겹 층화 교차 검증 (5-Fold Stratified CV)**: 데이터를 5개로 균등 분할하여 모델당 5회의 평가를 반복하고 그 평균치를 산출하였습니다. 이는 모델이 가려진(Unknown) 데이터에 대해서도 동일한 수준의 성능을 출력할 수 있음을 입증합니다.
- 특성 무결성 (Feature Integrity)**: n_segment, total_length 및 백분위수(Percentile) 데이터의 메타데이터를 엄격히 관리하는 Pandas 기반 파이프라인을 구축하여, 모델이 통계적 노이즈가 아닌 실질적 생물학적 특징을 학습하도록 설계하였습니다.
- 재현성 보증**: 벤치마크 기록 시점과 동일한 데이터 로딩 순서 및 전처리 로직을 파이프라인에 내재화하여, 결과의 변동성을 최소화하고 신뢰도를 극대화하였습니다.

3. 기술적 의견 및 향후 제언 (Technical Insights)

- GBDT 모델의 탁월성**: 현재의 고차원 유전체 정식 데이터(Tabular Data) 구조에서는 복잡한 신경망(MLP/CNN) 보다 그래디언트 부스팅(GBDT) 계열 모델이 학습 효율과 해석성 면에서 월등히 유리함을 확인하였습니다.
- 특성 활용 전략**: 1cM 임계값 적용은 정보 손실을 최소화하는 전략으로, 데이터의 미묘한 차이를 감지해야 하는 친족 분류 작업에 매우 적합한 선택이었습니다.
- 성능 고도화 여력**: 향후 정확도 97% 이상을 목표로 할 경우, 1cM 이하의 초미세 세그먼트 분석이나 특정 염색체 별 가중치 도입 등의 추가 연구를 제안합니다. 다만, 이는 데이터 노이즈 필터링 고도화 작업이 병행되어야 합니다.

4. 종합 결론 (Final Conclusion)

현재 구축된 **LightGBM** 기반 5겹 평가 파이프라인은 96.50%라는 최상위권 성능과 논리적 신뢰성을 동시에 확보하였습니다. 본 모델은 데이터 내의 풍부한 정보량을 효과적으로 추출하고 있으며, 현재의 결과는 실무 적용 및 향후 고도화 연구의 강력한 기반이 될 것으로 판단됩니다.