□ 머신러닝 모델 분석 및 정의서

NSMC 한국어 영화 리뷰 감정분석 프로젝트

🗐 문서 개요

문서 목적

본 문서는 NSMC 한국어 영화 리뷰 감정분석 프로젝트에서 사용된 머신러닝 알고리즘들의 이론적 배경, 실제 구 현, 성능 분석을 체계적으로 정리한 기술 문서입니다.

대상 독자

- 초급~중급 데이터 사이언티스트
- 머신러닝 학습자
- 텍스트 분류 프로젝트 수행자
- 프로젝트 리뷰어 및 평가자

문서 구성

- 각 알고리즘별 이론적 배경
- NSMC 데이터에서의 실제 성능
- 장단점 및 적용 시나리오
- 실무 활용 가이드라인

1. 나이브 베이즈 (Naive Bayes)

1.1 알고리즘 정의

기본 원리: 베이즈 정리를 기반으로 한 확률적 분류 알고리즘

P(긍정|단어들) = P(단어들|긍정) × P(긍정) / P(단어들)

"나이브(Naive)" 가정: 모든 특성(단어)이 서로 독립적이라고 가정

- 실제로는 단어들이 연관되어 있지만, 독립적이라고 "순진하게" 가정
- 이 가정 덕분에 계산이 매우 빠르고 효율적

1.2 NSMC 프로젝트 적용

사용한 변형: MultinomialNB (다항 나이브 베이즈)

- 텍스트 분류에 특화된 버전
- 단어의 출현 빈도를 고려

하이퍼파라미터:

python

MultinomialNB(alpha=1.0)

• (alpha=1.0): 라플라스 스무딩 (0확률 방지)

1.3 성능 결과

지표	수치	비고	
교차 검증 정확도	79.35% ± 0.27%	매우 안정적	
검증 데이터 정확도	79.25%	최고 성능	
학습 시간	0.014초	초고속	
F1-Score (긍정)	0.78	균형잡힌 성능	
F1-Score (부정)	0.80	우수한 부정 분류	
∢		▶	

1.4 핵심 발견사항

강력한 감정 지표 단어들:

긍정 신호:

- "굿굿": 99.2% 긍정 확률

- "재미있었어요": 99.2% 긍정 확률

- "ㅎㅎ": 84.2% 긍정 확률

부정 신호:

- "1점도": 98.7% 부정 확률

- "최악의 영화": 98.8% 부정 확률

- "???": 70.3% 부정 확률

한국어 특성 발견:

- 이모티콘이 강력한 감정 신호 역할
- 구어체 표현의 높은 구분력
- 짧은 텍스트에서도 명확한 감정 표현

1.5 장단점 분석

✓ 장점:

- 1. **초고속 학습**: 0.014초 (실시간 적용 가능)
- 2. **높은 성능**: 79.25% 정확도
- 3. 완벽한 해석성: 단어별 확률 직접 확인
- 4. 메모리 효율성: 적은 메모리 사용
- 5. 텍스트 특화: 문서 분류에 최적화
- 6. **안정성**: 일관된 성능 (표준편차 0.27%)

⚠ 한계:

- 1. 독립성 가정: 단어 간 관계 무시
- 2. 단순한 패턴: 복잡한 언어 패턴 놓칠 수 있음
- 3. 맥락 이해 부족: "기분좋게 틀었는데 잡쳤다" 같은 복합 표현 어려움

1.6 적용 시나리오

◎ 최적 적용 상황:

- 실시간 감정분석 (채팅, 댓글 모니터링)
- 대용량 배치 처리
- 프로토타입 개발
- 리소스 제약 환경
- 명확한 감정 표현이 있는 짧은 텍스트

★ 부적합 상황:

- 복잡한 문학 작품 분석
- 아이러니, 반어법이 많은 텍스트
- 긴 문서의 미묘한 감정 분석

☑ 2. 로지스틱 회귀 (Logistic Regression)

2.1 알고리즘 정의

기본 원리: 선형 조합 + 시그모이드 함수를 통한 확률 예측

$$z = w_1 \times x_1 + w_2 \times x_2 + ... + w_n \times x_n + b$$

P(긍정) = 1 / (1 + e^(-z))

핵심 특징:

• 각 단어에 가중치(weight) 부여

- 모든 단어의 가중치를 선형 조합
- 시그모이드 함수로 0-1 확률 변환

2.2 NSMC 프로젝트 적용

사용한 설정:

```
python

LogisticRegression(
C=1.0, # 정규화 강도
max_iter=1000, # 최대 반복 횟수
solver='liblinear' # 희소 행렬 최적화
)
```

최적화 과정:

- 6번 반복으로 정상 수렴
- 994번의 여유분 확보
- 안정적인 해법 도달

2.3 성능 결과

지표	수치	비교
교차 검증 정확도	79.13% ± 0.28%	나이브 베이즈와 유사
검증 데이터 정확도	79.17%	0.08%p 차이
학습 시간	0.604초	43배 더 오래
수렴 반복 횟수	6번	빠른 수렴
F1-Score (긍정)	0.78	동등한 성능

2.4 가중치 분석 결과

최고 긍정 기여 단어들:

```
1. "여운이": +5.26 (최고 긍정 지표)
2. "최고": +4.97 (직접적 긍정)
3. "최고의": +4.93
4. "수작": +4.45 (한국어 특유 표현)
5. "명작": +3.93 (작품성 평가)
```

최고 부정 기여 단어들:

- 1. "최악의": -7.00 (최강 부정 지표)
- 2. "최악": -5.78
- 3. "아깝다": -5.72 (한국어 특유 후회 표현)
- 4. "쓰레기": -5.49 5. "재미없다": -5.03

가중치 특성:

- 총 10,000개 특성 중 99.1%가 의미있는 가중치
- 가중치 범위: -7.00 ~ +5.26
- 표준편차: 1.01 (적절한 분산)

2.5 모델 해석성 분석

나이브 베이즈와의 공통점:

- "여운이", "재미있어요", "좋았어요" (긍정)
- "최악", "아깝다", "재미없다" (부정)

로지스틱 회귀만의 관점:

- 작품성 평가 단어에 높은 가중치 ("명작", "수작")
- 감정적 반응 단어 중시 ("가슴이")
- 더 세밀한 가중치 차별화

2.6 예측 차이 사례

개선된 맥락 이해:

예시: "초라해진 그에게 연민을 느낀다" 실제: 부정

나이브 베이즈: 긍정 (58.3% 확신) **※** 로지스틱 회귀: 부정 (51.0% 확신) ✓

분석: "초라해진"의 부정적 맥락을 가중치 조합으로 더 잘 포착

2.7 장단점 분석

✓ 장점:

- 1. **완벽한 해석성**: 가중치로 단어별 기여도 명확히 확인
- 2. 균형잡힌 성능: 나이브 베이즈와 동등한 79.17% 정확도
- 3. 빠른 수렴: 6번 반복으로 최적해 도달

- 4. 확률 출력: 예측 신뢰도 확인 가능
- 5. 맥락 이해: 단어 조합의 전체적 패턴 고려
- 6. **정규화**: 과적합 방지 내장

⚠ 한계:

- 1. **학습 시간**: 나이브 베이즈 대비 43배 오래 (0.604초)
- 2. 선형 모델: 비선형 패턴 포착 한계
- 3. **특성 의존성**: TF-IDF 품질에 성능 좌우

2.8 적용 시나리오

◎ 최적 적용 상황:

- 비즈니스 인사이트 도출이 중요한 경우
- 모델 설명이 필요한 상황 (규제 업종)
- 단어별 영향력 분석이 필요한 경우
- 성능과 해석성의 균형점이 필요한 경우

✓ 추천 상황:

- 마케팅 키워드 분석
- 고객 피드백 분석
- 브랜드 모니터링
- 설명 가능한 AI가 필요한 경우

1 3. SVM (Support Vector Machine)

3.1 알고리즘 정의

기본 원리: 마진을 최대화하는 최적 결정 경계 탐색

목표: 긍정과 부정을 나누는 고차원 초평면 찾기 조건: 마진(여백)을 최대화하여 일반화 성능 확보

핵심 개념:

- 서포트 벡터: 결정 경계에 가장 가까운 중요한 데이터들
- 마진: 결정 경계와 서포트 벡터 사이의 안전거리
- 커널: 비선형 패턴을 처리하기 위한 변환 함수

3.2 NSMC 프로젝트 적용

사용한 설정:

```
python

SVC(
kernel='linear', # 선형 커널
C=1.0, # 정규화 강도
probability=True # 확률 출력
)
```

선형 커널 선택 이유:

- TF-IDF 고차원 데이터 (10,000차원)에서 선형도 충분히 강력
- 빠른 학습 및 예측
- 해석 가능한 가중치 확인

3.3 성능 결과

지표	수치	문제점
교차 검증 정확도	78.77% ± 0.27%	가장 낮은 성능
검증 데이터 정확도	78.77%	나이브 베이즈보다 0.48%p 낮음
학습 시간	2,558초 (42분)	비현실적 느림
서포트 벡터 비율	52.7%	과도한 복잡성
효율성 점수	0.03	최저 효율성
◀	•	•

3.4 서포트 벡터 분석

서포트 벡터 특성:

• 전체 데이터의 52.7%가 서포트 벡터

• 부정 클래스: 약 25,000개

• 긍정 클래스: 약 25,000개

• 총 50,000개 이상의 서포트 벡터

문제점 분석:

이상적인 SVM: 10-20% 서포트 벡터

현재 상황: 50%+ 서포트 벡터

의미: 데이터가 선형적으로 잘 분리되지 않음

결과: SVM의 장점인 "적은 서포트 벡터로 효율적 분류" 무효화

3.5 NSMC 데이터에서 SVM 실패 원인

1. 데이터 특성 부적합:

- 평균 35.2자의 매우 짧은 텍스트
- 명확한 감정 구분 단어들
- 복잡한 패턴보다는 단순한 선형 분리가 효과적

2. 과도한 복잡성:

- 10,000차원 공간에서 9,999차원 초평면 탐색
- 필요 이상의 정교한 경계선 추구
- 단순한 문제를 복잡하게 접근

3. 계산 복잡도:

- O(n²~n³)의 시간 복잡도
- 119,963개 샘플에서 과도한 계산량
- 실시간 적용 불가능

3.6 제외 결정 사유

💥 성능 관점:

- 78.77% vs 나이브 베이즈 79.25% (0.48%p 낮음)
- 가장 낮은 정확도
- 성능 향상 효과 없음

💥 효율성 관점:

- 42분 학습 시간 (나이브 베이즈의 21,672배)
- 비현실적인 학습 속도
- 프로덕션 환경 부적합

💢 복잡성 관점:

- 52.7% 서포트 벡터 (과도한 복잡성)
- 메모리 사용량 과다
- 유지보수 어려움

3.7 SVM이 유리한 상황 (참고)

- 복잡한 비선형 패턴이 있는 데이터
- 중간 규모의 데이터셋 (1,000~10,000개)

- 고차원에서 선형 분리가 어려운 경우
- 정교한 결정 경계가 필요한 상황

예시:

- 의료 진단 데이터
- 이미지 분류 (작은 데이터셋)
- 바이오인포매틱스
- 복잡한 패턴의 시계열 데이터

📊 4. 모델 비교 종합 분석

4.1 성능 비교 요약

모델	정확도	학습시간	효율성*	해석성	복잡도	권장도
나이브 베이즈	79.25%	0.014초	5,661	☆☆☆☆	낮음	>
로지스틱 회귀	79.17%	0.604초	131	$^{\diamond}$	중간	S
SVM	78.77%	2,558초	0.03	☆☆	높음	×
◀	•	•	-	•	•	>

^{*}효율성 = 정확도 / 학습시간 × 1000

4.2 NSMC 데이터 특성에 따른 모델 적합성

데이터 특성:

√ 평균 길이: 35.2자 (매우 짧음)

✓ 평균 단어 수: 7.6개 (단순함)

✓ 감정 구분: 명확한 지표 단어

√ 언어적 특성: 구어체, 이모티콘 활용

✓ 패턴: 선형적 분리 가능

모델별 적합성 평가:

나이브 베이즈: 🥎 🏠 🥎 🥎

- 짧은 텍스트에 최적화
- 독립성 가정이 크게 위배되지 않음
- 명확한 감정 단어들을 확률로 잘 포착

로지스틱 회귀: 🔷 🔷 🔷 🔷

- 단어 조합의 선형 관계 포착
- 가중치로 명확한 해석 제공

• 적당한 복잡도로 균형점 제공

SVM: 🗙 🏠

- 단순한 데이터에 과도한 복잡성
- 필요 이상의 정교한 경계선 탐색
- 계산 비용 대비 성능 향상 없음

4.3 실무 관점 종합 평가

프로덕션 환경 고려사항:

요구사항	나이브 베이즈	로지스틱 회귀	SVM
실시간 처리	✓ 최적	☑ 가능	💢 불가능
배치 처리	✓ 최적	✓ 우수	💢 비효율
해석 필요성	✓ 우수	✓ 최고	⚠ 어려움
리소스 제약	✓ 최적	☑ 양호	💢 과다
유지보수	✓ 쉬움	☑ 쉬움	💢 복잡
4	•	•	>

◎ 5. 실무 활용 가이드라인

5.1 상황별 모델 선택 가이드

☑ 실시간 감정분석 (나이브 베이즈 추천)

적용 사례:

- 실시간 댓글 모니터링
- 소셜미디어 감정 추이 분석
- 고객 서비스 채팅 분석
- 뉴스 댓글 실시간 분류

선택 이유:

- √ 0.014초 초고속 처리
- √ 79.25% 높은 정확도
- √ 메모리 효율적
- √ 확장성 우수

○ 비즈니스 인사이트 도출 (로지스틱 회귀 추천)

★ SVM 비추천 상황

✓ 동등한 성능 (79.17%)

NSMC 같은 데이터에서는 비추천:

- 짧고 단순한 텍스트
- 명확한 감정 구분 단어
- 실시간 처리 요구
- 리소스 제약 환경
- 해석성이 중요한 경우

5.2 성능 최적화 전략

나이브 베이즈 최적화:

```
python
# 추천 설정
MultinomialNB(
alpha=1.0, # 기본값 유지
fit_prior=True # 클래스 분포 고려
)
# 성능 향상 팁
- TF-IDF 전처리 품질 향상
- 한국어 불용어 정제
- N-gram 조합 실험 (1,2)
```

로지스틱 회귀 최적화:

python

5.3 모델 해석 및 모니터링

나이브 베이즈 해석 방법:

```
python
# 단어별 확률 확인
feature_log_prob = model.feature_log_prob_
positive_words = feature_names[np.argsort(feature_log_prob[1])[-20:]]
# 예측 신뢰도 분석
probabilities = model.predict_proba(X_test)
confidence_scores = np.max(probabilities, axis=1)
```

로지스틱 회귀 해석 방법:

5.4 배포 및 운영 고려사항

모델 저장 및 로딩:

python

```
# 모델 저장
import joblib
joblib.dump(model, 'nsmc_naive_bayes_v1.pkl')
joblib.dump(tfidf_vectorizer, 'nsmc_tfidf_v1.pkl')

# 모델 로딩
model = joblib.load('nsmc_naive_bayes_v1.pkl')
vectorizer = joblib.load('nsmc_tfidf_v1.pkl')
```

성능 모니터링:

```
python
# 성능 저하 감지

def monitor_model_performance(model, new_data, threshold=0.75):
    accuracy = model.score(new_data.X, new_data.y)
    if accuracy < threshold:
        alert_performance_degradation(accuracy)
    return accuracy
```

모델 업데이트 전략:

- 월간 성능 평가
- 새로운 데이터 패턴 모니터링
- A/B 테스트를 통한 모델 개선
- 점진적 업데이트 (Blue-Green 배포)

뎚 6. 결론 및 권장사항

6.1 핵심 결론

🛂 최적 모델: 나이브 베이즈

- NSMC 데이터 특성에 가장 적합
- 최고 성능 (79.25%) + 초고속 (0.014초)
- 실무 환경에서 즉시 적용 가능

🥑 대안 모델: 로지스틱 회귀

- 해석성이 중요한 경우 최적 선택
- 나이브 베이즈와 동등한 성능 (79.17%)
- 비즈니스 인사이트 도출에 유리

💥 비추천: SVM

- NSMC 데이터에는 과도한 복잡성
- 성능 대비 비효율적 자원 사용
- 다른 유형의 데이터에서 재고려

6.2 프로젝트 학습 성과

이론적 성과:

- 1. "No Free Lunch" 실증: 복잡한 모델이 항상 좋지 않음
- 2. 데이터 기반 선택: 데이터 특성에 맞는 알고리즘 중요성
- 3. 효율성 평가: 성능과 자원 사용의 균형점

실무적 성과:

- 1. 체계적 비교: 교차 검증을 통한 안정적 평가
- 2. 해석 능력: 모델 결과의 비즈니스 활용 방안
- 3. 배포 준비: 프로덕션 환경 고려사항 정리

6.3 향후 발전 방향

단기 개선 계획:

- 하이퍼파라미터 세밀 튜닝
- 앙상블 모델 시도 (나이브 베이즈 + 로지스틱 회귀)
- 한국어 특화 전처리 강화

장기 발전 계획:

- 딥러닝 모델 도전 (LSTM, CNN)
- Transformer 모델 적용 (KoBERT)
- 85%+ 정확도 달성 목표

6.4 최종 권장사항

프로덕션 배포 시:

- 1. **1순위**: 나이브 베이즈 (실시간 처리)
- 2. 2순위: 로지스틱 회귀 (해석성 중요시)
- 3. **병행 운영**: A/B 테스트로 실제 환경 검증

학습 및 연구 목적:

• 3개 모델 모두 구현하여 알고리즘 특성 이해

- 각 모델의 장단점 체험으로 모델 선택 안목 향상
- 실무 지향적 사고 프로세스 습득

🗐 부록

A. 주요 하이퍼파라미터 설명

나이브 베이즈:

- (alpha): 라플라스 스무딩 강도 (기본값: 1.0)
- (fit_prior): 클래스 확률 학습 여부 (기본값: True)

로지스틱 회귀:

- C: 정규화 강도 (작을수록 강한 정규화)
- max_iter]: 최대 반복 횟수
- (solver): 최적화 알고리즘 ('liblinear' 추천)

SVM:

- (C): 마진 vs 오분류 트레이드오프
- (kernel): 커널 함수 ('linear', 'rbf', 'poly')
- probability : 확률 출력 여부

B. 성능 지표 해석 가이드

정확도 (Accuracy): 전체 예측 중 올바른 예측 비율 정밀도 (Precision): 긍정 예측 중 실제 긍정 비율 재현율 (Recall): 실제 긍정 중 올바르게 예측한 비율 F1-Score: 정밀도와 재현율의 조화평균

C. 추가 참고 자료

- NSMC 데이터셋: <u>https://github.com/e9t/nsmc</u>
- scikit-learn 문서: <u>https://scikit-learn.org/</u>
- 한국어 자연어처리: KoNLPy, KoBERT
- 텍스트 전처리: TF-IDF, Word2Vec

문서 버전: v1.0 작성일: 2025년 1월 작성자: NSMC 감정분석 프로젝트팀 검토일: 2025년 1월