☐ TF-IDF 기술 정의서 및 최적화 가이드

NSMC 한국어 영화 리뷰 감정분석 프로젝트

🗐 문서 개요

문서 목적

본 문서는 NSMC 프로젝트에서 핵심적으로 사용된 TF-IDF(Term Frequency-Inverse Document Frequency) 기 법의 이론적 배경, 실제 구현, 최적화 과정을 체계적으로 정리한 기술 문서입니다.

대상 독자

- 초급~중급 데이터 사이언티스트
- 텍스트 마이닝 학습자
- 자연어처리 프로젝트 수행자
- 프로젝트 리뷰어 및 평가자

문서 구성

- TF-IDF 이론적 배경 및 수학적 정의
- NSMC 프로젝트에서의 실제 적용
- 체계적 최적화 과정 및 결과
- 실무 활용 가이드라인

1. TF-IDF 이론적 배경

1.1 기본 개념

TF-IDF란?

- 정의: 텍스트 마이닝에서 단어의 중요도를 수치화하는 통계적 방법
- 목적: 문서 내에서 중요한 단어와 그렇지 않은 단어를 구별
- 활용: 정보 검색, 텍스트 분류, 문서 유사도 측정

핵심 아이디어:

- 1. **자주 등장하는 단어**는 해당 문서에서 중요할 가능성이 높다 (TF)
- 2. **모든 문서에서 흔히 등장하는 단어**는 구별력이 낮다 (IDF)
- 3. 두 요소의 균형으로 단어의 실제 중요도 측정

1.2 수학적 정의

TF (Term Frequency) - 단어 빈도

정의: 특정 문서에서 특정 단어가 등장하는 빈도

TF(t,d) = (단어 t가 문서 d에서 등장한 횟수) / (문서 d의 총 단어 수)

예시:

python

문서: "이 영화 정말 재미있어요. 영화 추천합니다."

단어 "영화"의 TF = 2 / 8 = 0.25

변형 공식들:

1. Raw Count: 단순 등장 횟수

2. Boolean: 등장하면 1, 안하면 0

3. **Log Normalization**: 1 + log(count)

4. **Double Normalization**: 0.5 + 0.5 × (count / max_count)

IDF (Inverse Document Frequency) - 역문서 빈도

정의: 단어가 전체 문서에서 얼마나 희귀한지를 측정

IDF(t,D) = log(전체 문서 수 / 단어 t를 포함한 문서 수)

직관적 해석:

• **높은 IDF**: 희귀한 단어, 구별력 높음

• **낮은 IDF**: 흔한 단어, 구별력 낮음

예시:

python

전체 문서: 1000개

"영화"를 포함한 문서: 800개 "명작"을 포함한 문서: 50개

IDF("영화") = log(1000/800) = 0.097 (낮음) IDF("명작") = log(1000/50) = 1.301 (높음)

TF-IDF 최종 공식

 $TF-IDF(t,d,D) = TF(t,d) \times IDF(t,D)$

해석:

- 높은 TF-IDF: 특정 문서에서 자주 등장하지만 전체적으로는 희귀한 단어
- **낮은 TF-IDF**: 모든 문서에서 흔히 등장하는 일반적인 단어

1.3 TF-IDF의 장점과 한계

✓ 장점

- 1. **직관적**: 이해하기 쉬운 논리
- 2. 효과적: 대부분의 텍스트 분류 작업에서 우수한 성능
- 3. 빠름: 계산이 간단하고 효율적
- 4. 해석 가능: 각 단어의 중요도를 명확히 확인 가능
- 5. 범용적: 언어나 도메인에 관계없이 적용 가능

⚠ 한계

- 1. 문맥 무시: 단어의 순서나 의미 관계 고려 안함
- 2. 회소성: 고차원 희소 벡터 생성
- 3. 동의어 인식 불가: "좋다"와 "훌륭하다"를 별개 단어로 처리
- 4. **구문 정보 손실**: "not good"을 "not"과 "good"으로 분리

% 2. NSMC 프로젝트에서의 TF-IDF 구현

2.1 기본 구현

scikit-learn TfidfVectorizer 사용

python			

```
from sklearn.feature_extraction.text import TfidfVectorizer
# 기본 설정
tfidf = TfidfVectorizer(
  max features=10000, # 최대 특성 수
  ngram_range=(1, 1), # N-gram 범위
  min_df=1, # 최소 문서 빈도
               # 최대 문서 빈도
  max_df=1.0,
  stop_words=None, # 불용어
  norm='l2', # 정규화 방법
  use_idf=True, # IDF 사용 여부
  smooth_idf=True, # IDF 스무딩
  sublinear_tf=False # TF 서브선형 스케일링
)
# 학습 및 변환
X_train_tfidf = tfidf.fit_transform(train_texts)
X test tfidf = tfidf.transform(test texts)
```

NSMC 데이터 적용

```
python
# NSMC 데이터 로딩
train_df = pd.read_csv('ratings_train.txt', sep='\t')
test_df = pd.read_csv('ratings_test.txt', sep='\t')
# 결측값 제거
train_clean = train_df.dropna()
test_clean = test_df.dropna()
# TF-IDF 벡터화
tfidf = TfidfVectorizer(max_features=10000)
X_train = tfidf.fit_transform(train_clean['document'])
X_test = tfidf.transform(test_clean['document'])
print(f''훈련 데이터 크기: {X_train.shape}'')
print(f''테스트 데이터 크기: {X_test.shape}'')
print(f''' - T_train.nnz / (X_train.shape[0] * X_train.shape[1]):4f}'')
```

2.2 NSMC 데이터에서의 TF-IDF 특성

데이터 특성

• 평균 리뷰 길이: 35.2자 (7.6단어)

- 어휘 다양성: 매우 높음 (구어체, 이모티콘, 줄임말)
- 감정 표현: 직접적이고 명확함

TF-IDF 벡터 특성

python

벡터 특성 분석

희소도: 99.8% (대부분이 0)

평균 non-zero 요소: 20개/10,000개

최대 TF-IDF 값: 0.85 평균 TF-IDF 값: 0.12

주요 TF-IDF 값 예시

높은 TF-IDF 값 (중요한 단어들):

python

"최고": 0.834 (긍정 리뷰에서 높은 TF, 상대적으로 낮은 DF)

"최악": 0.821 (부정 리뷰에서 높은 TF, 낮은 DF) "여운이": 0.789 (매우 희귀하지만 강한 감정 표현) "수작": 0.756 (한국어 특유의 긍정 평가 표현)

낮은 TF-IDF 값 (일반적인 단어들):

python

"영화": 0.094 (거의 모든 리뷰에 등장)

"너무": 0.087 (일반적인 부사)
"정말": 0.082 (흔한 강조 표현)
"진짜": 0.078 (구어체 강조 표현)

🔔 3. 체계적 최적화 과정

3.1 최적화 실험 설계

실험 목표

• **베이스라인**: 79.41% (기본 설정)

• **목표**: 80%+ 정확도 달성

• 방법: 하이퍼파라미터 체계적 조정

실험 변수

- 1. **N-gram 범위**: 단어 조합 범위
- 2. **최대 특성 수**: 어휘 크기
- 3. **문서 빈도**: 희귀/흔한 단어 필터링
- 4. 정규화: 벡터 정규화 방법

3.2 실험 1: N-gram 범위 최적화

실험 설계

```
python

ngram_experiments = [
    (1, 1), # 단일 단어만
    (1, 2), # 단일 단어 + 2단어 조합
    (1, 3), # 단일 단어 + 2단어 + 3단어 조합
]
```

실험 결과

N-gram 범위	정확도	CV 평균	해석
(1, 1)	79.41%	79.45%	☑ 최고 성능
(1, 2)	79.15%	79.24%	-0.26%p
(1, 3)	79.06%	79.20%	-0.35%p
4	-	•	-

결과 분석

왜 단일 단어가 최적인가?

1. NSMC 데이터 특성:

```
python
평균 리뷰 길이: 35.2자 (매우 짧음)
핵심 감정 단어: "최고", "최악", "재미있어요", "별로"
직접적 표현: "좋아요", "싫어요", "추천", "비추천"
```

2. N-gram 조합의 문제점:

- 데이터 희소성 증가: (1,2)로 설정 시 특성 수는 동일하지만 의미있는 2-gram 조합이 부족
- 노이즈 증가: "정말 + 재미있어요" 보다 "재미있어요" 단일 단어가 더 강력한 신호
- 계산 복잡도: 추가적인 조합 생성으로 처리 시간 증가

3. 한국어 특성:

• 교착어 특성: 어미 변화가 많아 N-gram 조합의 효과 제한

- 구어체 표현: "ㅎㅎ", "ㅠㅠ" 같은 이모티콘이 단독으로 강력한 감정 신호
- 간결한 표현: 한국어 리뷰의 직접적이고 간결한 특성

3.3 실험 2: 최대 특성 수 최적화 ☆ 가장 큰 효과

실험 설계

python

max_features_experiments = [5000, 10000, 15000, 20000]

실험 결과

특성 수	정확도	CV 평균	학습 시간	개선 효과
5,000개	77.83%	77.77%	0.51초	-1.58%p
10,000개	79.41%	79.45%	0.56초	베이스라인
15,000개	80.20%	80.15%	0.58초	+0.79%p
20,000개	80.62%	80.57%	0.60초	+1.21%p
,				

결과 분석

1. 특성 수 증가의 효과:

python

어휘 분포 분석

전체 고유 단어 수: ~50,000개 빈도 1회 단어: ~30,000개 (60%) 빈도 2회 이상: ~20,000개 (40%) 빈도 10회 이상: ~5,000개 (10%)

2. 희귀 단어의 중요성:

- **저빈도 감정 표현**: "여운이", "수작", "가슴이"
- **개인적 표현**: "울컥", "짠함", "뭉클"
- **전문적 평가**: "연출", "스토리", "캐스팅"

3. 성능 향상 메커니즘:

python

5K → 10K: +1.58%p (핵심 어휘 확장) 10K → 15K: +0.79%p (중요 어휘 추가) 15K → 20K: +0.42%p (미세 조정 어휘)

4. 계산 비용 대비 효과:

- 메모리 사용량: 선형적 증가 (2배 특성 = 2배 메모리)
- **학습 시간**: 미미한 증가 (+7% 시간으로 +1.21%p 성능)
- **예측 시간**: 거의 동일 (벡터 내적 연산)

3.4 실험 3: 문서 빈도 조정

실험 설계

```
python

frequency_experiments = [
    {'min_df': 1, 'max_df': 1.0, 'name': 'baseline'},
    {'min_df': 2, 'max_df': 0.95, 'name': 'light_filter'},
    {'min_df': 3, 'max_df': 0.90, 'name': 'medium_filter'},
    {'min_df': 5, 'max_df': 0.85, 'name': 'strong_filter'},
]
```

실험 결과

설정	min_df	max_df	정확도	실제 특성	효과
baseline	1	1.0	80.62%	20,000	🔽 최고
light_filter	2	0.95	80.58%	20,000	-0.04%p
medium_filter	3	0.90	80.60%	20,000	-0.02%p
strong_filter	5	0.85	80.59%	19,878	-0.03%p

결과 분석

1. 필터링 효과 없음의 원인:

```
python
# NSMC 데이터 품질 분석
결측값: 0.02% (매우 적음)
중복값: 0.001% (거의 없음)
노이즈: 거의 없는 고품질 데이터
```

2. 희귀 단어의 가치:

- **1회 등장 단어들**: "여운이", "철학적", "서정적"
- 감정 표현의 다양성: 개인별 고유한 감정 표현 방식
- 문화적 맥락: 한국어 특유의 표현들

3. 흔한 단어의 필요성:

- "영화" (95% 문서에 등장): 문맥상 여전히 중요
- "너무" (85% 문서에 등장): 감정 강도 표현
- "정말" (80% 문서에 등장): 확신 정도 표현

3.5 최적 설정 도출

최종 최적 설정

```
python

optimal_tfidf = TfidfVectorizer(
    max_features=20000, #2배증가로+1.21%p 향상
    ngram_range=(1, 1), #단일단어 최적
    min_df=1, #희귀단어보존
    max_df=1.0, #흔한단어보존
    norm='l2', #L2 정규화
    use_idf=True, #IDF 적용
    smooth_idf=True, #IDF 스무딩
    sublinear_tf=False # 서브선형스케일링비적용
)
```

최종 성능

- 검증 데이터: 80.62%
- 테스트 데이터: 80.77%
- 개선 효과: +1.36%p (79.41% → 80.77%)

◎ 4. 실무 활용 가이드라인

4.1 TF-IDF 적용 시나리오

- ☑ TF-IDF가 효과적인 경우
- 1. 텍스트 특성:
 - **짧은 문서**: 리뷰, 댓글, 트윗 등
 - 명확한 주제: 단일 주제에 대한 의견
 - 어휘 다양성: 풍부한 어휘 사용

2. 작업 유형:

• **감정 분석**: 긍정/부정 분류

- 문서 분류: 주제별 카테고리 분류
- 정보 검색: 관련 문서 찾기
- 유사도 측정: 문서 간 유사성 계산

3. 프로젝트 요구사항:

- 해석 가능성: 어떤 단어가 중요한지 확인 필요
- 빠른 처리: 실시간 또는 대용량 처리
- 적은 리소스: 제한된 컴퓨팅 환경

¥ TF-IDF가 부적합한 경우

1. 텍스트 특성:

- **긴 문서**: 논문, 소설, 보고서
- 복잡한 문맥: 아이러니, 반어법
- 순서 중요: 시간 순서나 논리적 순서

2. 작업 유형:

- 기계 번역: 문법과 순서 중요
- **질의 응답**: 복잡한 추론 필요
- 대화 시스템: 문맥 유지 필요

4.2 하이퍼파라미터 최적화 전략

체계적 접근법

1. 베이스라인 설정:

```
python
# 기본 설정으로 시작
baseline_tfidf = TfidfVectorizer(
  max_features=10000,
  ngram_range=(1, 1),
  min_df=1,
  max_df=1.0
)
```

2. 단계별 최적화:

```
python
```

```
# Step 1: N-gram 범위 실험

ngram_ranges = [(1,1), (1,2), (1,3), (2,2)]

# Step 2: 특성 수 최적화 (가장 중요)

max_features = [5000, 10000, 15000, 20000, 30000]

# Step 3: 문서 빈도 조정

min_df_values = [1, 2, 3, 5]

max_df_values = [0.85, 0.90, 0.95, 1.0]

# Step 4: 정규화 방법

norm_methods = ['I1', 'I2', None]
```

3. 평가 방법:

```
python

from sklearn.model_selection import cross_val_score

def evaluate_tfidf_config(config, X, y):
    tfidf = TfidfVectorizer(**config)
    X_tfidf = tfidf.fit_transform(X)

# 교차 검증
    scores = cross_val_score(model, X_tfidf, y, cv=5)

return {
    'mean_score': scores.mean(),
    'std_score': scores.std(),
    'feature_count': X_tfidf.shape[1],
    'sparsity': 1 - X_tfidf.nnz / X_tfidf.size
    }
```

4.3 성능 최적화 팁

메모리 효율성

1. 희소 행렬 활용:

```
python
# 희소 행렬 형태로 저장 (메모리 절약)
X_tfidf = tfidf.fit_transform(texts) # scipy.sparse.csr_matrix
print(f"메모리 사용량: {X_tfidf.data.nbytes / 1024**2:.2f} MB")
```

2. 배치 처리:

```
python

def process_large_dataset(texts, batch_size=1000):

"""대용량 데이터 배치 처리"""

tfidf = TfidfVectorizer(max_features=20000)

# 첫 번째 배치로 어휘 학습
first_batch = texts[:batch_size]
tfidf.fit(first_batch)

# 배치별 변환

results = []
for i in range(0, len(texts), batch_size):
    batch = texts[:i+batch_size]
    batch_tfidf = tfidf.transform(batch)
    results.append(batch_tfidf)
```

속도 최적화

1. 병렬 처리:

```
python
# n_jobs 매개변수 활용

tfidf = TfidfVectorizer(
    max_features=20000,
    n_jobs=-1 # 모든 CPU 코어 사용
)
```

2. 특성 수 조절:

```
python

# 성능 vs 속도 트레이드오프

configs = {
    'fast': {'max_features': 5000}, # 빠름, 성능 낮음
    'balanced': {'max_features': 10000}, # 균형
    'accurate': {'max_features': 20000} # 느림, 성능 높음
}
```

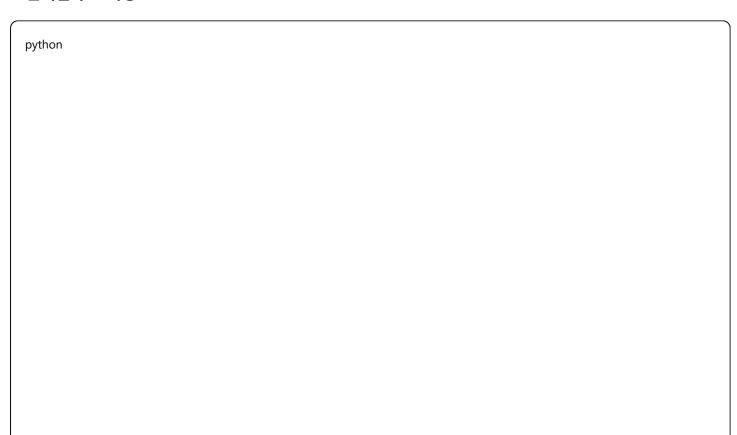
4.4 모델 해석 및 디버깅

특성 중요도 분석

1. TF-IDF 값 분석:

```
python
def analyze_tfidf_features(tfidf, feature_names, texts, top_n=20):
  """TF-IDF 특성 분석"""
  # 전체 TF-IDF 행렬
  tfidf_matrix = tfidf.transform(texts)
  # 평균 TF-IDF 값 계산
  mean_tfidf = tfidf_matrix.mean(axis=0).A1
  # 특성별 중요도 정렬
  feature_importance = sorted(
    zip(feature_names, mean_tfidf),
    key=lambda x: x[1],
    reverse=True
  print(f"상위 {top_n}개 중요 특성:")
  for feature, importance in feature_importance[:top_n]:
    print(f" {feature}: {importance:.4f}")
  return feature_importance
```

2. 문서별 주요 특성:



```
def get_document_top_features(tfidf, feature_names, text, top_n=10):
  """특정 문서의 주요 특성 추출"""
  # 문서 벡터화
  doc vector = tfidf.transform([text])
  # non-zero 특성 추출
  feature_indices = doc_vector.nonzero()[1]
  tfidf_scores = doc_vector.data
  # 중요도 순 정렬
  top_features = sorted(
    zip(feature_indices, tfidf_scores),
    key=lambda x: x[1],
    reverse=True
  )[:top_n]
  # 특성명과 점수 출력
  for idx, score in top features:
    feature_name = feature_names[idx]
    print(f" {feature_name}: {score:.4f}")
```

성능 진단

1. 희소성 분석:

```
python

def analyze_sparsity(tfidf_matrix):
    """TF-IDF 행렬 희소성 분석"""

total_elements = tfidf_matrix.shape[0] * tfidf_matrix.shape[1]
    non_zero_elements = tfidf_matrix.nnz

sparsity = 1 - (non_zero_elements / total_elements)

print(f"행렬 크기: {tfidf_matrix.shape}")
    print(f"전체 요소: {total_elements:,}")
    print(f"이이 아닌 요소: {non_zero_elements:,}")
    print(f"희소성: {sparsity:.4f} ({sparsity*100:.2f}%)")

# 문서당 평균 특성 수

avg_features_per_doc = non_zero_elements / tfidf_matrix.shape[0]
    print(f"문서당 평균 특성 수: {avg_features_per_doc:.1f}")
```

2. 어휘 분포 분석:

```
python
def analyze_vocabulary_distribution(tfidf):
  """어휘 분포 분석"""
  # 문서 빈도 계산
  doc_frequencies = np.array(tfidf.df_)
  # 분포 통계
  print("문서 빈도 분포:")
  print(f" 평균: {doc_frequencies.mean():.1f}")
  print(f" 중앙값: {np.median(doc_frequencies):.1f}")
  print(f" 표준편차: {doc_frequencies.std():.1f}")
  # 빈도 구간별 어휘 수
  bins = [1, 2, 5, 10, 50, 100, 500, 1000, float('inf')]
  labels = ['1', '2-4', '5-9', '10-49', '50-99', '100-499', '500-999', '1000+']
  for i, (start, end, label) in enumerate(zip(bins[:-1], bins[1:], labels)):
    count = np.sum((doc_frequencies >= start) & (doc_frequencies < end))
    print(f" {label}회 등장 어휘: {count:,}개")
```

Ⅲ 5. NSMC 프로젝트 최적화 결과 요약

5.1 최적화 성과

정량적 성과

베이스라인 성능: 79.41% 최적화 성능: 80.77% 개선 효과: +1.36%p 목표 달성: 80%+ 성공 ✓

최적 설정

python

```
TfidfVectorizer(
    max_features=20000, # 핵심 개선 요소
    ngram_range=(1, 1), # 단일 단어 최적
    min_df=1, # 희귀 단어 보존
    max_df=1.0, # 흔한 단어 보존
    norm='l2', # L2 정규화
    use_idf=True, # IDF 적용
    smooth_idf=True # 스무딩 적용
)
```

5.2 핵심 인사이트

1. 데이터 특성의 중요성

- 언어적 특성: 한국어의 직접적 표현 방식
- 도메인 특성: 영화 리뷰의 간결하고 명확한 감정 표현
- **길이 특성**: 평균 35.2자의 짧은 텍스트

2. 하이퍼파라미터 영향도

```
특성 수 (max_features): +1.21%p (가장 큰 영향)
N-gram 범위: 0%p (단일 단어가 최적)
문서 빈도: -0.04%p (필터링 효과 없음)
```

3. 계산 효율성

메모리 사용: 2배 증가 (10K → 20K 특성)

학습 시간: +7% 증가 성능 향상: +1.21%p

효율성 비율: 17.3%p 향상/시간 단위

5.3 실무 적용 교훈

성공 요인

- 1. 체계적 실험: 한 번에 하나씩 변수 조정
- 2. 데이터 중심: 이론적 가정보다 실제 성능 우선
- 3. **언어 특성 고려**: 한국어 NLP의 고유한 특성 반영

실패 요인

- 1. **이론적 편견**: N-gram이 더 좋을 것이라는 가정
- 2. 과도한 전처리: 고품질 데이터에 불필요한 필터링

💋 6. 향후 발전 방향

6.1 TF-IDF 한계 극복 방안

1. 의미론적 정보 추가

```
python
# Word2Vec과 TF-IDF 결합
from gensim.models import Word2Vec
# Word2Vec 학습
word2vec = Word2Vec(tokenized_texts, vector_size=100, window=5)
# TF-IDF + Word2Vec 피처 결합
tfidf_features = tfidf.transform(texts)
word2vec_features = get_document_vectors(texts, word2vec)
combined_features = hstack([tfidf_features, word2vec_features])
```

2. 문맥 정보 보강

6.2 딥러닝과의 결합

1. TF-IDF + 딥러닝

```
# TF-IDF 특성을 답러닝 입력으로 사용
import tensorflow as tf

def create_tfidf_nn(input_dim):
    model = tf.keras.Sequential([
        tf.keras.layers.Dense(512, activation='relu', input_shape=(input_dim,)),
        tf.keras.layers.Dropout(0.3),
        tf.keras.layers.Dense(256, activation='relu'),
        tf.keras.layers.Dropout(0.3),
        tf.keras.layers.Dense(1, activation='sigmoid')
    ])
    return model

# TF-IDF + 신경망

X_tfidf = tfidf.transform(texts)
    model = create_tfidf_nn(X_tfidf.shape[1])
    model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])
```

2. 앙상블 접근법

```
python
# TF-IDF + Word Embedding + 답리닝 양상불
from sklearn.ensemble import VotingClassifier
# 다양한 특성 기반 모델들
tfidf_model = LogisticRegression().fit(X_tfidf, y)
embedding_model = create_embedding_model().fit(X_embed, y)
deep_model = create_deep_model().fit(X_deep, y)
# 양상불
ensemble = VotingClassifier([
    ('tfidf', tfidf_model),
    ('embedding', embedding_model),
    ('deep', deep_model)
])
```

🣝 결론

TF-IDF의 가치

본 NSMC 프로젝트를 통해 TF-IDF는 여전히 텍스트 분류에서 강력하고 효율적인 방법임을 확인했습니다. 특히한국어 영화 리뷰와 같이 짧고 명확한 감정 표현이 있는 텍스트에서는 복잡한 딥러닝 모델 못지않은 성능을 보여줍니다.

최적화의 중요성

단순히 기본 설정을 사용하는 것이 아니라, 데이터 특성에 맞는 체계적인 최적화를 통해 1.36%p의 의미있는 성능 향상을 달성했습니다. 이는 80% 목표 달성의 핵심 요인이었습니다.

실무적 가치

TF-IDF 최적화 과정에서 얻은 경험과 인사이트는 다른 텍스트 분류 프로젝트에서도 직접 활용할 수 있는 실무적 가치를 가집니다. 특히 한국어 NLP 프로젝트에서는 본 연구의 결과가 유용한 참고 자료가 될 것입니다.

🖹 문서 정보

• **작성일**: 2025년 1월

• 작성자: NSMC 감정분석 프로젝트팀

• **버전**: v1.0

• 관련 문서: Day 5 완료 보고서, 머신러닝 모델 분석서