



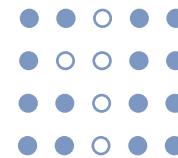
# R – Project

## 도서 분야 예측 시스템

아이티윌 빅데이터 31기  
김동환

# INDEX

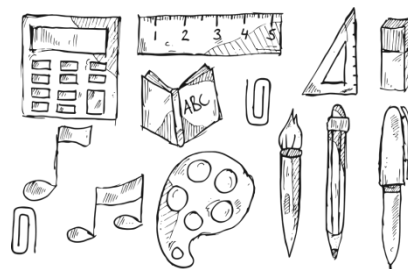
Enjoy your stylish business and campus life with BIZCAM



1. 배경 및 목표



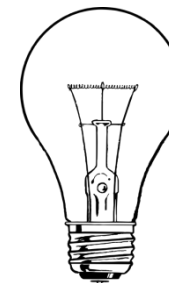
2. 구현 순서



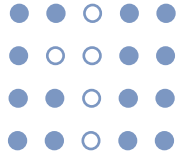
3. 주요 내용



4. 결론

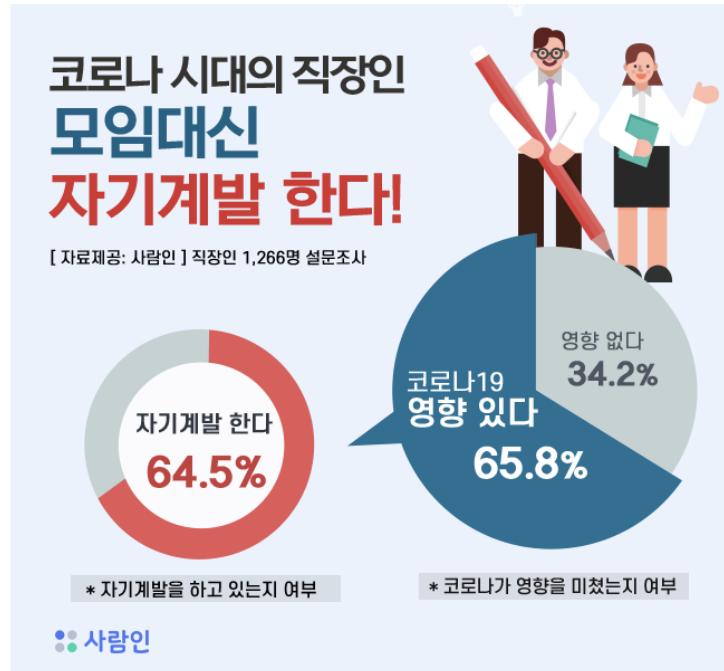


5. 방향성



## 코로나 19시대의 직장인, 모임 대신 '자기계발'한다!

등록일 : 2021.05.10 조회수 : 1,959



Copyrights © 사람인.



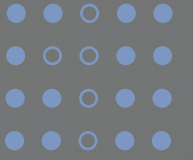
독서는 우리를 즐겁게 하고 독특한 경험을 제공하기 때문에 많은 사람들이 좋아하는 취미이다. 좋은 책은 독자에게 환상의 분위기를 조성할 수 있는 힘을 가지고 있으며, 다른 나라 및 세계를 방문하고, 다른 배경을 가진 사람들을 만나고, 그들의 운명을 발견하도록 이끌어 주는 매체이다.

독서는 정신을 위해 할 수 있는 가장 건강한 일 중 하나라는 것을 누구나 알고 있으며 많은 사람들이 휴식이나 정보를 위해 책을 읽지만 그 외의 독서가 우리에게 과학적으로 주는 모든 이점이란 무엇인지 알아가 보자.

Coyrights © parisunni 코코파.  
All right reserved 블로그

# 목표

Enjoy your stylish business and campus life with BIZCAM

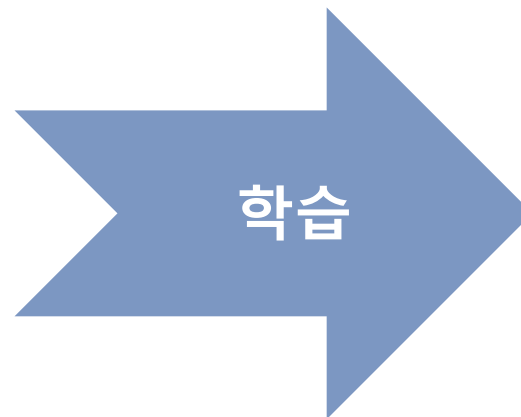
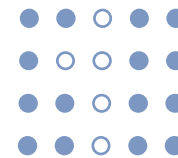


각 분야별 문서의 제목과 소개 내용을 추출 및 학습하고  
분석하여 새로운 문서 텍스트에 대한 분야 예측



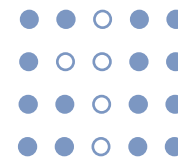
# 구현 순서

Enjoy your stylish business and campus life with BIZCAM



# 주요 라이브러리 / 함수

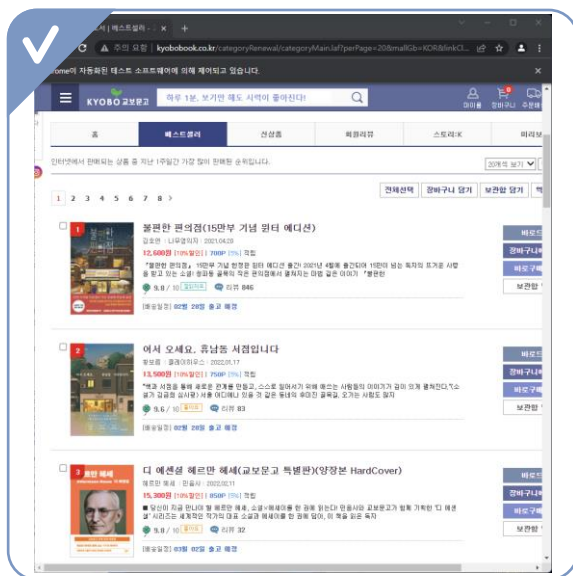
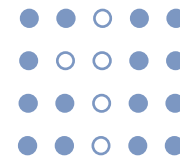
Enjoy your stylish business and campus life with BIZCAM



- **rvest** # 웹페이지에서 html태그를 추출하고 태그의 다양한 옵션들을 사용할 수 있는 라이브러리
  - ✓ `Read_html`, `html_nodes`, `html_text`
- **RSelenium** # 동적인 웹페이지의 태그를 추출할 수 있는 라이브러리
  - ✓ `remoteDriver$(open, navigate, findElement, getPageSource, goBack)`
- **tm** # VCorpus(말뭉치, documentTermMatrix 만드는 작업을 수행하는 라이브러리
  - ✓ `VCorpus`, `DocumentTermMatrix`, `Terms`, `inspect`, `findFreqTerms`
- **stringr** # 단어 추출, 변경, 제거, 공백제거 등 전처리작업에 주로 사용
  - ✓ `str_replace_all`, `str_match_all`,
- **RcppMeCab** # 일본어, 한국어 품사 태깅 라이브러리
  - ✓ `pos`
- **e1071** # naiveBayes모델을 사용할 수 있는 라이브러리
  - ✓ `naiveBayes`
- **gmodels** # 두 개의 질적 자료 간의 관련성을 교차표로 나타내줄 수 있는 함수 `CrossTable`을 사용할 수 있는 라이브러리
  - ✓ `CrossTable`

# 데이터 수집

Enjoy your stylish business and campus life with BIZCAM



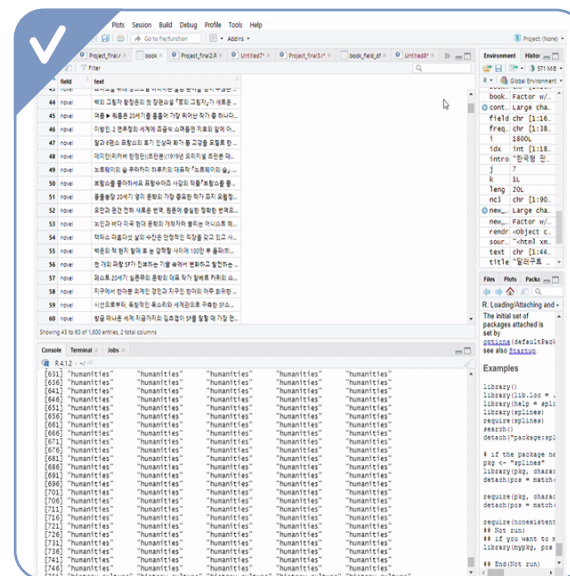
1

각 장르별 베스트셀러의  
url 추출



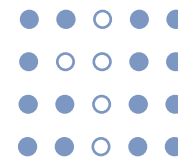
2

각 책의 url에 들어가서  
제목과 소개글을 추출



3

정답라벨, 추출한 텍스트  
데이터프레임으로 저장



```
#말뭉치로 변형작업
book_corpus <- tm::VCorpus(VectorSource(book$text))

#태깅, 전처리 작업
#토큰으로 만들 함수
mecab_words <- function(doc){
  # 태깅
  tagging <- RcppMeCab::pos(base::enc2utf8(as.character(doc)))

  # 필요한 품사 NNG(일반명사), NNP(고유명사), XR(어근 : 실질적으로 의미를 지니는 형태소)
  word_noun <- str_match_all(tagging, '(\w+)/NNG|NNP|XR|SL')[[1]][,2]

  # 전처리
  word_noun <- str_replace_all(word_noun, 'c', '씨언어')
  word_noun <- str_replace_all(word_noun, 'R', '알프로그래밍')
  word_noun <- str_replace_all(word_noun, 'NFT', '엔에프티')
  word_noun <- str_replace_all(word_noun, 'ETF', '이티에프')
  word_noun <- str_replace_all(word_noun, 'LH', '한국토지주택공사')
```

sd): 0





```
<<DocumentTermMatrix (documents: 1800, terms: 14330)>>
Non-/sparse entries: 98369/25695631
Sparsity           : 100%
Maximal term length: 15
Weighting           : binary (bin)
Sample             :
      Terms
docs  내용 독자 문제 사람 세계 시작 이야기 이해 자신 저자
1103  1    0    1    0    0    0    0    0    1    0
1124  0    0    1    1    1    1    1    1    1    1
1159  0    0    1    1    1    0    0    0    0    1
1182  0    0    1    1    0    1    1    0    1    0
1417  1    1    1    1    1    0    1    1    1    0
3     0    1    0    0    1    1    1    1    0    1
36    0    1    1    1    1    1    1    1    0    1
5     0    1    0    0    1    0    1    1    0    1
504   0    0    1    1    1    1    1    1    0    1
867   0    0    1    1    1    0    1    1    0    0
```

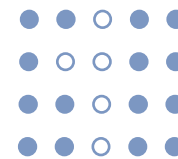


```
[idx [1] 1 1 1 1 1 1 1 1 1 1 2 2 1 1 1 1 1 1 1 1 1 1 1 1 2  
[51] 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 2 1 1 1 1 1 1  
101] 1 1 1 1 2 1 1 2 1 1 1 2 1 1 2 2 1 1 1 1 2 1 2 1 1 1 1 1 2  
151] 1 1 1 1 1 1 1 2 2 1 2 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2  
201] 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1  
251] 2 1 1 1 1 1 1 1 1 1 2 1 1 1 2 1 2 1 1 1 1 1 1 2 1 1 1 1  
301] 1 1 1 1 1 1 2 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1  
351] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 2 1 1 2 1 1 1 1  
401] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1  
451] 1 1 1 1 1 1 2 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 2 1 1 1  
501] 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
551] 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
601] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
651] 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
701] 1 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1  
751] 1 1 1 1 1 2 1 1 1 1 2 1 1 1 1 1 2 1 1 1 1 2 1 1 1 1 1  
801] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 1 1 2 1 2 1  
851] 1 1 1 1 1 1 1 2 1 1 1 1 1 2 1 1 1 2 1 1 1 2 1 1 1 1 1  
901] 1 1 1 2 1 1 1 1 1 1 1 1 1 2 1 1 1 1 1 1 2 2 1 1 1 1 1  
951] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 1 1 2 1 1 1
```

```
reached getOption("max.print") -- omitted 800 entries ]
```

# 학습 / 테스트 데이터 분류

Enjoy your stylish business and campus life with BIZCAM



## 학습데이터

```
<<DocumentTermMatrix (documents: 1624, terms: 14330)>>
Non-/sparse entries: 89223/23182697
Sparsity : 100%
Maximal term length: 15
Weighting : binary (bin)
Sample :
  Terms
Docs  내용 독자 문제 사람 세계 시작 이야기 이해 자신 저자
1103  1    0    1    0    0    0    0    0    1    0    0
1124  0    0    1    1    1    1    1    1    1    1    1
1159  0    0    1    1    1    1    0    0    0    1    1
1182  0    0    1    1    1    0    1    1    0    1    0
1417  1    1    1    1    1    0    1    1    1    0    1
3     0    1    0    0    1    1    1    1    0    1    0
36    0    1    1    1    1    1    1    1    0    1    0
5     0    1    0    0    1    0    1    1    0    1    0
504   0    0    1    1    1    1    1    1    0    1    1
867   0    0    1    1    1    0    1    1    0    0    1
```

학습데이터 개수

```
> nrow(book_dtm_train)
[1] 1624
```

학습데이터 정답라벨 개수

```
> length(book_dtm_train_labels)
[1] 1624
```

## 테스트데이터

```
<<DocumentTermMatrix (documents: 176, terms: 14330)>>
Non-/sparse entries: 9146/2512934
Sparsity : 100%
Maximal term length: 15
Weighting : binary (bin)
Sample :
  Terms
Docs  내용 문제 사람 세계 시작 이야기 이해 자신 저자 필요
1095  1    1    0    0    0    0    0    0    0    0    1
1136  0    1    0    1    0    0    0    1    0    0    0
1199  0    1    0    1    1    0    0    0    0    0    0
1292  1    1    1    0    0    0    0    1    0    1    0
1602  0    0    0    0    0    0    0    1    0    0    0
1662  1    1    0    0    0    0    0    1    0    1    1
235   1    0    0    1    1    1    1    0    0    0    0
335   0    0    1    0    1    0    0    0    1    1    1
346   0    0    0    1    0    1    1    1    0    0    0
68    0    1    0    1    0    1    1    0    1    0    0
```

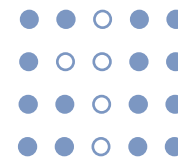
테스트데이터 개수    테스트데이터 정답라벨 개수

```
> nrow(book_dtm_test)
[1] 176
```

```
> length(book_dtm_test_labels)
[1] 176
```

# 학습 / 테스트 데이터 분류

Enjoy your stylish business and campus life with BIZCAM



## apply

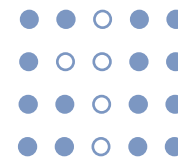
```
convert_counts <- function(x){  
  x <- ifelse(x>0, 'YES', 'NO')  
}
```

```
book_train <- apply(book_dtm_train, MARGIN=2, convert_counts)  
book_test <- apply(book_dtm_test, MARGIN=2, convert_counts)
```

가가	가가미	가감	가계	가격	가계
"NO"	"NO"	"NO"	"NO"	"NO"	"NO"
가급	가짜이	가나안	가난	가네코	가늠
"NO"	"NO"	"NO"	"NO"	"NO"	"NO"
가득	가득성	가드	가드너	가든	가디언
"NO"	"NO"	"NO"	"NO"	"NO"	"NO"
로등	가로세로	가로수	가루	가르	가르시아
"NO"	"NO"	"NO"	"NO"	"NO"	"NO"
가명	가문	가름	가미	가미야	가미오카
"NO"	"NO"	"NO"	"NO"	"NO"	"NO"
플라	가비지	가사	가산점	가상	가상현실
"NO"	"NO"	"NO"	"NO"	"NO"	"NO"
가속	가속도	가수	가스	가스관	가스실
"NO"	"NO"	"NO"	"NO"	"NO"	"NO"
가시	가식	가신	가액	가야	가예
"NO"	"NO"	"NO"	"NO"	"NO"	"NO"
가운	가운데	가을	가의	가이	가이드
"NO"	"NO"	"NO"	"NO"	"NO"	"NO"

# naiveBayes, predict

Enjoy your stylish business and campus life with BIZCAM



## • naiveBayes

```
book_classifier <- naiveBayes(book_train,book_dtm_train_labels)
```

book_dtm_train_labels	특역사	
	NO	YES
art_popularculture	1.00000000	0.00000000
computer_it	1.00000000	0.00000000
economy	1.00000000	0.00000000
history_culture	0.98591549	0.01408451
humanities	1.00000000	0.00000000
novel	1.00000000	0.00000000
poem_essay	1.00000000	0.00000000
politics_society	1.00000000	0.00000000
religion	1.00000000	0.00000000
selfdevelope	1.00000000	0.00000000
sience	1.00000000	0.00000000
technology_engineering	1.00000000	0.00000000

book_dtm_train_labels	힌두교	
	NO	YES
art_popularculture	1.00000000	0.00000000
computer_it	1.00000000	0.00000000
economy	1.00000000	0.00000000
history_culture	0.98496241	0.01503759
humanities	1.00000000	0.00000000
novel	1.00000000	0.00000000
poem_essay	1.00000000	0.00000000
politics_society	1.00000000	0.00000000
religion	0.99270073	0.00729927
selfdevelope	1.00000000	0.00000000
sience	1.00000000	0.00000000
technology_engineering	1.00000000	0.00000000

book_dtm_train_labels	힌트	
	NO	YES
art_popularculture	1.00000000	0.00000000
computer_it	1.00000000	0.00000000
economy	0.992537313	0.007462687
history_culture	1.00000000	0.00000000
humanities	1.00000000	0.00000000
novel	1.00000000	0.00000000
poem_essay	0.992366412	0.007633588
politics_society	1.00000000	0.00000000
religion	1.00000000	0.00000000
selfdevelope	0.978571429	0.021428571
sience	1.00000000	0.00000000
technology_engineering	1.00000000	0.00000000

## • predict

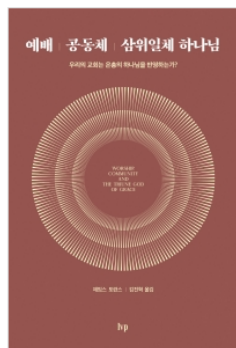
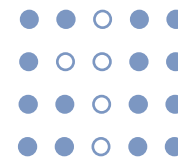
```
book_test_predict <- predict(book_classifier,book_test)
```

```
> book_test_predict
[1] poem_essay      novel            novel            politics_society
[5] novel           poem_essay      novel            novel
[9] novel           history_culture novel            novel
[13] novel           poem_essay      poem_essay       poem_essay
[17] art_popularculture poem_essay      poem_essay       poem_essay
[21] poem_essay      poem_essay      poem_essay       history_culture
[25] poem_essay      humanities      poem_essay       poem_essay
[29] selfdevelope    economy        economy          economy
[33] economy         economy        selfdevelope     economy
[37] politics_society economy        economy          economy
[41] economy         politics_society economy          humanities
[45] economy         economy        selfdevelope     selfdevelope
[49] poem_essay      religion        selfdevelope     humanities
[53] selfdevelope    humanities      selfdevelope     selfdevelope
[57] selfdevelope    poem_essay      selfdevelope     selfdevelope
[61] religion        politics_society selfdevelope     politics_society
[65] humanities      humanities      selfdevelope     humanities
```



# 추가 테스트 데이터

Enjoy your stylish business and campus life with BIZCAM



무료배송 | 소독공제

**[출간예정] 예배, 공동체, 삼위일체 하나님**  
가?

제임스 토런스 지음 | 김진혁 옮김 | IVP | 2022년 03월 04일 출간

정가: 12,000원

판매가: **10,800원** [10%↓ 1,200원 할인]

혜택: [기본적립] 600원 적립 [5% 적립]

[추가적립] 5만원 이상 구매 시 **2,000원** 추가적립 [안내](#)

[회원혜택] 회원 등급 별, 3만원 이상 구매 시 **2~4%** 추가적립 [안내](#)

## 책소개

이 책이 속한 분야

종교 > 기독교(개신교) > 기독교일반 > 기독교일반

예배, 성례, 교제의 참된 모습은 무엇인가?

삼위일체 하나님이 교회에 불러넣으신 참생명이 드러나다!

“예배 현장에서 교리가 나왔다. 그러나 바른 교리는 바른 예배의 초석이 된다.

이 책은 우리로 하여금 바른 예배에 눈뜨게 할 것이다!”

문화당 고려신학대학원 예배학 교수

제임스 토런스가 삼위일체 교리를 토대로 예배 신학에 관해 강의한 1994년 디즈버리 강연에 기초하여 출간한 책이다. 예배에 대한 새로운 이해가 필요한 시대, 우리의 예전은 무엇에 기초하는가. 제임스 토런스는 전통 교리와 히브리서에 대한 깊은 탄탄한 이해를 바탕으로, 예배가 은총의 삼위일체 하나님의 존재와 성자 그리스도 예수님이 행하신 바를 반영해야 함을 밝힌다. 그리고 은총의 선물로 주어진 예배에 우리가 성령으로 참여하게 된다는 복된 소식을 전한다. 이 책은 그 분량의 간결함에도 그리스도의 대제사장직 교리, 세례와 성찬에 대한 이해와 실행, 현대 페미니즘에 대한 반응까지 신학과 성례전 해석을 부족함 없이 다룬다. 이제 교회는 예수 그리스도 안에 있는 참된 중심으로 돌아가라는 이 책의 외침에 귀 기울여야 한다.

<http://www.kyobobook.co.kr/product/detailViewKor.laf?mallGb=KOR&ejkGb=KOR&linkClass=21030301&barcode=9788932819143>

corpus



Terms



DocumentTermMatrix

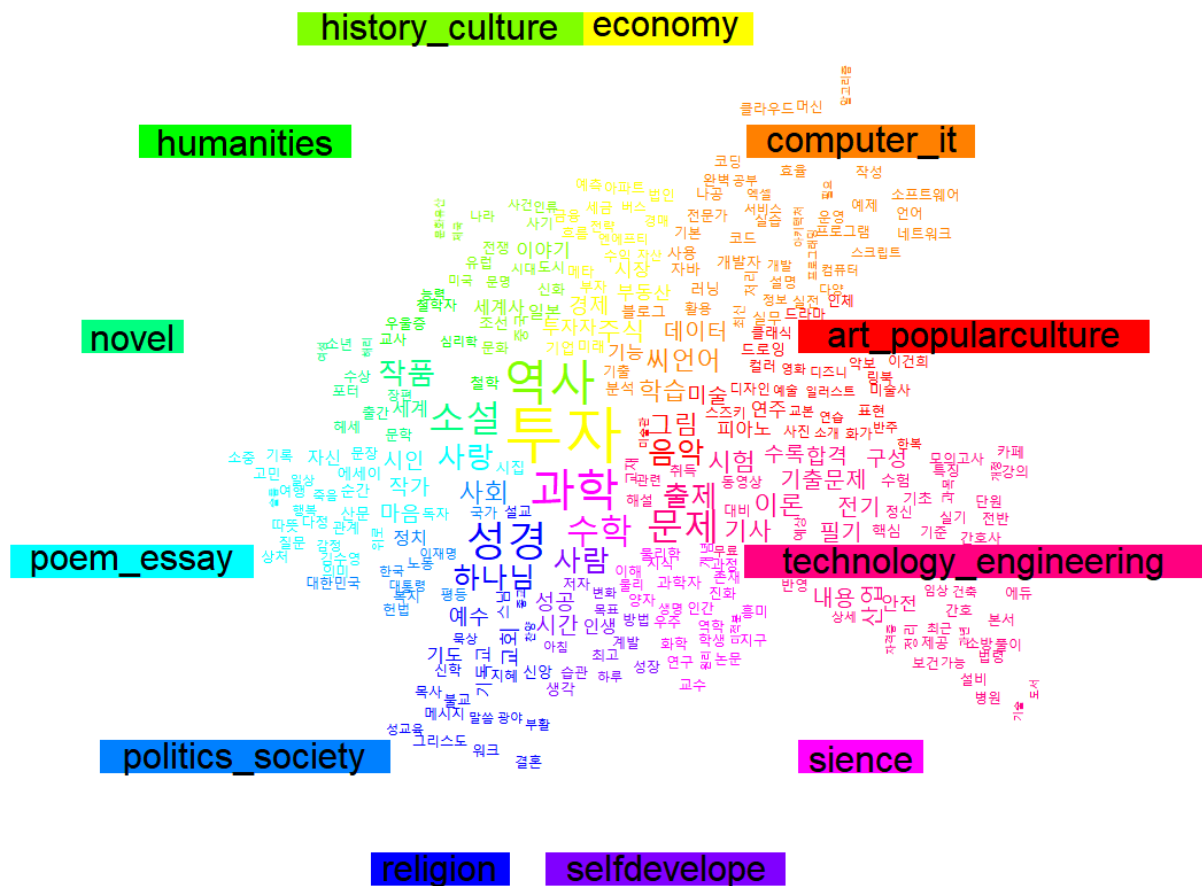


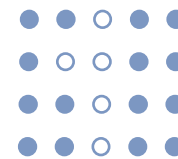
apply 이분화



```
> predict(book_classifier,t(new_dtm_test))
[1] religion
Levels: art_popularculture computer_it ecoi
~ |
```

— □ ×





1. 기존의 낮은 예측률 문제 해결
2. 사람들의 상황, 상태에 맞춤형 책을 추천해주는 시스템 구현
  - 평점, 후기, 소개, 목차, 제목을 분석 및 학습하여 text를 입력했을 때 text와 관련된 책 추천



감사합니다