

SeverusPT automated text clasification

JG, APP

1 Text mining and classification of papers

The main objective for this task is to use examples of papers that are selected (or not) for each of the phases of wildfire research, namely:

- Pre-fire
- During-fire
- Post-fire

Then a Random Forest model will be used to learn with these binary examples and classify all the remaining papers as selected or non-selected.

1.1 Count functions

Load libs and define ancillary functions

1.2 Data

Read the data and perform some basic preparation

Check colnames in the dataset

[1]	"...1"	"label"
[3]	"label_1"	"type"
[5]	"author"	"title"
[7]	"year"	"journal"
[9]	"volume"	"number"
[11]	"pages"	"doi"
[13]	"url"	"affiliations"
[15]	"abstract"	"author_keywords"
[17]	"correspondence_address"	"publisher"

```

[19] "issn" "language"
[21] "abbrev_source_title" "publication_stage"
[23] "source" "note"
[25] "keywords" "coden"
[27] "pmid" "isbn"
[29] "editor" "address"
[31] "affiliation" "earlyaccessdate"
[33] "eissn" "keywords_plus"
[35] "research_areas" "web_of_science_categories"
[37] "author_email" "orcid_numbers"
[39] "cited_references" "number_of_cited_references"
[41] "times_cited" "usage_count_last_180_days"
[43] "usage_count_since_2013" "journal_iso"
[45] "doc_delivery_number" "web_of_science_index"
[47] "unique_id" "month"
[49] "researcherid_numbers" "article_number"
[51] "funding_acknowledgement" "funding_text"
[53] "booktitle" "series"
[55] "book_author" "filename"
[57] "n_duplicates" "oa"
[59] "da" "citation"
[61] "screened_titles" "notes"
[63] "included"

```

Print class percentages:

Var1	Freq
excluded	48.32
selected	51.68

1.3 Training data

The training dataset is made by combining the following binary features,

- Top-30 most common words/terms from the **title** in *selected papers*;
- Top-30 most common words/terms from the **title** in *non-selected papers*;
- Top-30 most common words/terms from the **keywords** in *selected papers*;
- Top-30 most common words/terms from the **keywords** in *non-selected papers*;

1.3.1 Step 1 - count the most frequent terms in the title and keywords

Title top-30 terms

Count words in titles (after removing common/stop words) and select those terms that are most common in selected/accepted papers and non-selected/non-accepted papers.

These terms will be used as binary features in RF classification.

List of common terms in *titles* for *selected papers*:

included	word	n
1	fire	1039
1	forest	292
1	management	276
1	prescribed	178
1	risk	168
1	regimes	159
1	wildfire	159
1	weather	113
1	effects	106
1	forests	96
1	regime	95
1	fires	92
1	climate	79
1	fuel	66
1	model	55
1	change	54
1	prediction	54
1	vegetation	54
1	pine	52
1	carbon	51
1	soil	50
1	hazard	49
1	landscape	49
1	post	48
1	australia	47
1	implications	47
1	assessment	44
1	usa	44
1	mediterranean	41
1	response	41

List of common terms in *titles* for *non-selected papers*:

included	word	n
0	protection	296
0	prevention	91
0	hazards	83
0	study	66
0	treatment	61
0	safety	58
0	concrete	56
0	analysis	54
0	intumescent	53
0	steel	53
0	blight	50
0	coatings	50
0	strategies	42
0	composites	41
0	passive	40
0	performance	40
0	epoxy	37
0	evaluation	36
0	thermal	36
0	power	35
0	systems	35
0	properties	34
0	resistance	34
0	reducing	32
0	system	32
0	coal	30
0	coating	30
0	nuclear	30
0	storage	30
0	plants	27

Keywords top-30 terms

Now repeat the same process but now using keywords:

List of common terms in *keywords* for *selected papers*:

included	keyword	n
1	fire	1421
1	forest	896
1	management	712
1	fires	504
1	climate	405
1	wildfire	376
1	ecosystem	295
1	change	290
1	united	273
1	analysis	263
1	risk	253
1	environmental	233
1	burning	217
1	prescribed	217
1	forestry	213
1	assessment	191
1	carbon	181
1	vegetation	171
1	species	162
1	model	160
1	soil	155
1	weather	155
1	conservation	144
1	land	143
1	australia	142
1	modeling	139
1	population	137
1	pinus	126
1	ecology	122
1	north	119

List of common terms in *keywords* for *non-selected papers*:

included	keyword	n
0	protection	420
0	safety	201
0	heat	186
0	human	181
0	article	176
0	thermal	163

included	keyword	n
0	health	162
0	resistance	160
0	flame	152
0	coatings	143
0	hazards	142
0	combustion	138
0	temperature	137
0	smoke	132
0	concrete	129
0	humans	116
0	study	115
0	extinguishers	109
0	coal	107
0	materials	104
0	steel	102
0	male	101
0	systems	100
0	adult	99
0	female	94
0	hazard	92
0	gas	85
0	water	84
0	aged	83
0	performance	83

Abstract top-30 terms

Now repeat the same process but now using the abstract:

List of common terms in *abstracts* for *selected papers*:

included	abstract__	n
1	fire	6983
1	forest	1621
1	fires	1264
1	management	1215
1	species	1013
1	model	798

included	abstract_	n
1	risk	781
1	wildfire	764
1	vegetation	753
1	fuel	735
1	prescribed	680
1	climate	656
1	study	645
1	results	618
1	burned	613
1	data	601
1	forests	565
1	burning	563
1	weather	558
1	soil	554
1	effects	534
1	regimes	508
1	treatments	477
1	severity	474
1	change	452
1	land	445
1	conditions	431
1	increased	420
1	time	395
1	landscape	386
1	increase	372
1	regime	369
1	sites	359
1	cover	356
1	low	353
1	frequency	351
1	burn	350
1	models	345
1	ecological	341
1	season	330
1	carbon	328
1	spatial	327
1	post	322
1	wildfires	317
1	tree	316
1	scale	313
1	plant	300

included	abstract__	n
1	potential	300
1	structure	294
1	future	288

List of common terms in *abstracts* for *non-selected papers*:

included	abstract__	n
0	protection	534
0	temperature	470
0	safety	396
0	heat	384
0	thermal	377
0	analysis	340
0	flame	309
0	paper	292
0	rate	268
0	performance	250
0	system	246
0	method	239
0	resistance	233
0	steel	231
0	design	217
0	materials	217
0	properties	217
0	coatings	215
0	effect	215
0	smoke	214
0	elsevier	213
0	test	213
0	compared	212
0	research	210
0	concrete	202
0	combustion	198
0	hazard	198
0	release	195
0	water	192
0	blight	188
0	rights	187
0	assessment	186
0	methods	186

included	abstract_	n
0	control	185
0	developed	182
0	reserved	181
0	coating	178
0	systems	174
0	significant	169
0	coal	166
0	process	165
0	oxygen	160
0	surface	159
0	effective	158
0	intumescent	158
0	retardant	158
0	approach	156
0	gas	154
0	treatment	153
0	experimental	148

List bigrams in *abstracts* for all papers:

bigram	n
fire risk	474
fire regimes	424
fire protection	375
prescribed fire	366
fire management	341
fire regime	331
forest fire	308
climate change	268
post fire	255
fire weather	251
fire hazard	203
fire frequency	190
forest fires	165
fire resistance	160
fire blight	157
fire safety	150
heat release	137
prescribed burning	135
fire suppression	134

bigram	n
fire activity	132
fire occurrence	130
risk assessment	129
fire severity	127
fire danger	107
fire prone	104
flame retardant	104
weather conditions	99
dry season	98
fire spread	96
wildfire risk	92
release rate	91
fire prevention	90
forest management	90
fire behavior	88
severity fire	88
prescribed fires	86
species richness	79
extreme fire	74
fire hazards	72
land management	72
land cover	70
fuel treatments	69
fire events	67
ponderosa pine	65
wildland fire	65
fire season	64
fire behaviour	63
boreal forest	60
historical fire	60
plant species	59
fire intensity	58
intumescent coatings	58
management strategies	57
national park	57
north america	57
pre fire	57
air quality	56
real time	56
vegetation types	56
fire history	55

bigram	n
fire size	55
remote sensing	55
short term	55
burn severity	54
fire return	54
forest structure	54
heat flux	54
fuel treatment	52
management practices	52
wind speed	52
fuel loads	50
fuel moisture	50
future fire	50
landscape scale	50
community composition	49
low intensity	49
spontaneous combustion	49
fire exclusion	48
fire fighting	47
prescribed burns	47
spatially explicit	47
weather index	47
western united	47
fuel load	46
intumescent coating	46
biomass burning	45
mixed conifer	45
peak heat	45
risk management	45
thermal stability	45
nuclear power	44
fire treatments	43
flame retardancy	43
passive fire	43
char layer	42
cone calorimeter	42
fire propagation	42
fire scenarios	42
mechanical properties	42
fire ant	41
fire ants	41

bigram	n
lithium ion	41
mass loss	41
potential fire	41
urban interface	41
ecosystem services	40
fire effects	40
fire related	40
future climate	40
sierra nevada	40
vegetation structure	40
fire retardant	39
growing season	39
crown fire	38
forest ecosystems	38
fuel reduction	38
information system	38
national forest	38
reduce fire	38
thermal conductivity	38
land managers	37
soil moisture	37
fire extinguishing	36
machine learning	36
species composition	36
total heat	36
erwinia amylovora	35
finite element	35
active fire	34
climatic conditions	34
fire performance	34
fire tests	34
low severity	34
natural fire	34
neural network	34
protection systems	34
severity fires	34
steel structures	34
wildland urban	34
carbon storage	33
conifer forests	33
fire needle	33

bigram	n
fire risks	33
increased fire	33
fire alarm	32
fire managers	32
northern australia	32
plant communities	32
plant community	32
scale fire	32
water mist	32
wildfire management	32
canopy cover	31
carbon stocks	31
fire prediction	31
mixed severity	31
natural disturbance	31
relative humidity	31
season fires	31
tree mortality	31
geographic information	30
heat transfer	30
pine forests	30

Generate bigrams for the entire dataset

1.3.2 Step 2 - make features

Count the selected words in the title and keywords and arrange them as a binary grid:

Title-based features

Keyword-based features

Abstract-based features

1.3.3 Step3 - assemble all features and labels

Make the training dataset by combining everything:

1.4 Prediction dataset

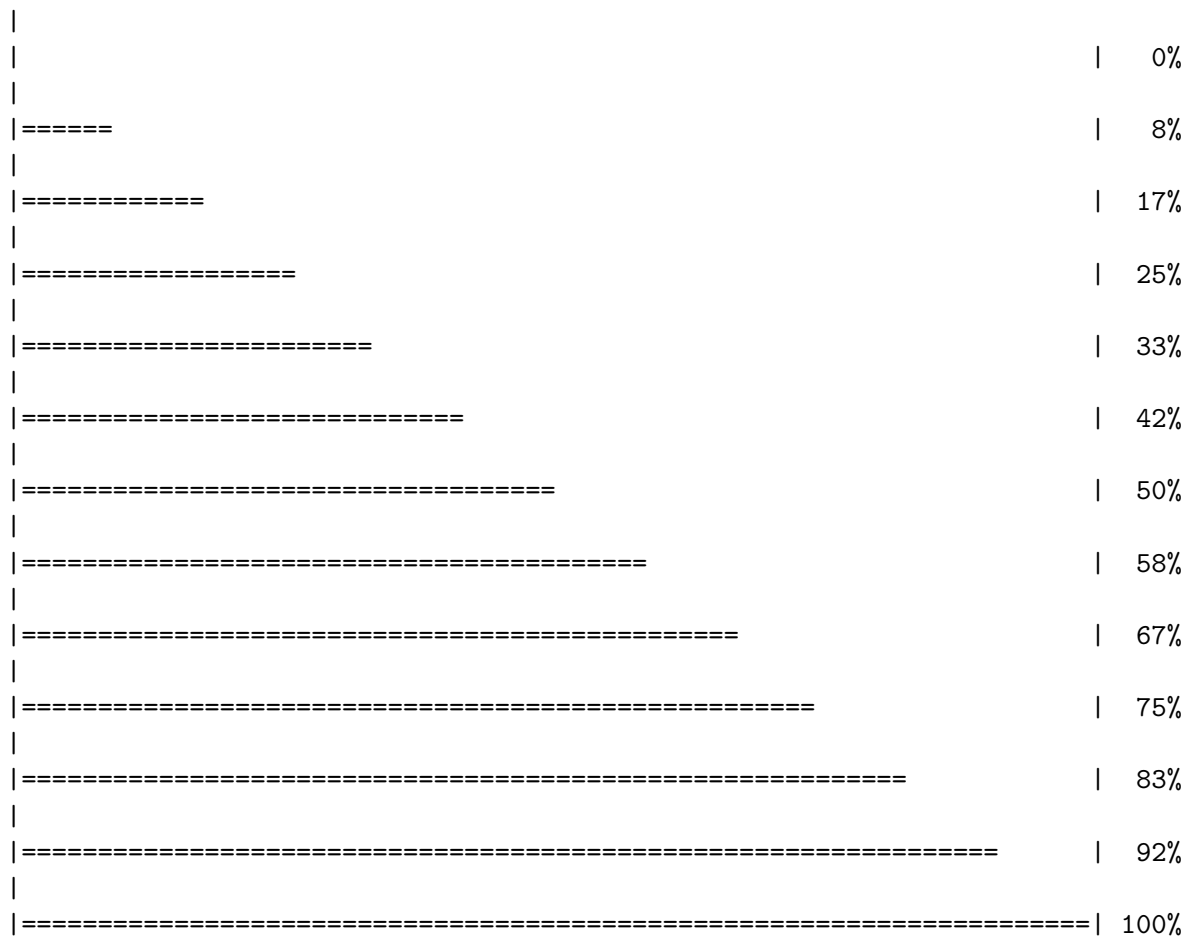
Step 1 - Evaluate the same features as before but now for the entire dataset of papers

Step 2 - Assemble the full prediction dataset:

1.5 Random Forest model development

1.5.1 Optimize RF hyperparameters

Make the classification model based on Random Forests:



best mtry value:

[1] 12

Call:

```
randomForest(x = train_tb %>% select(-label, -included), y = train_tb %>% pull(include
```

Type of random forest: classification

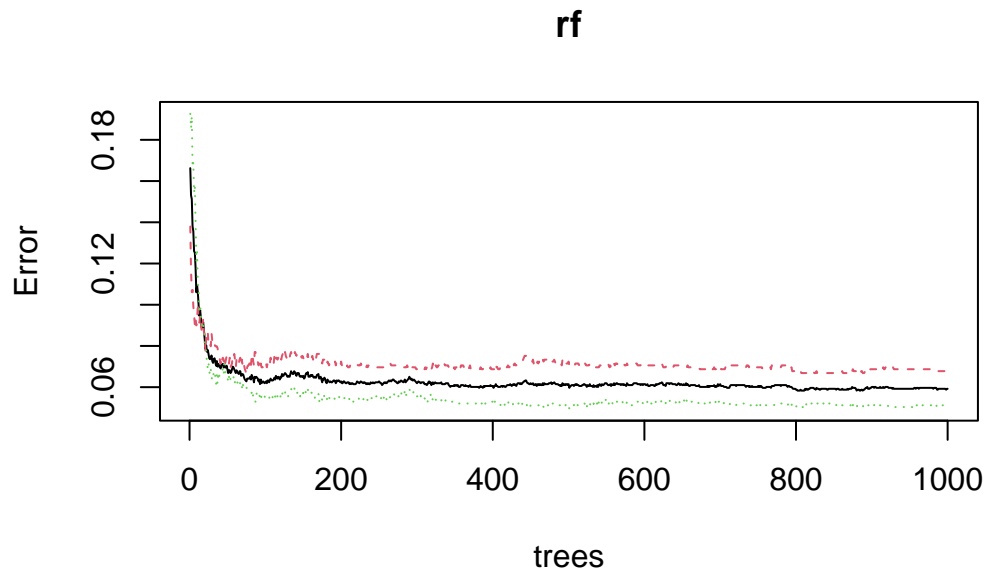
Number of trees: 1000

No. of variables tried at each split: 12

OOB estimate of error rate: 5.92%

Confusion matrix:

	0	1	class.error
0	1004	73	0.06778087
1	59	1093	0.05121528



1.5.2 Ten fold cross-validation

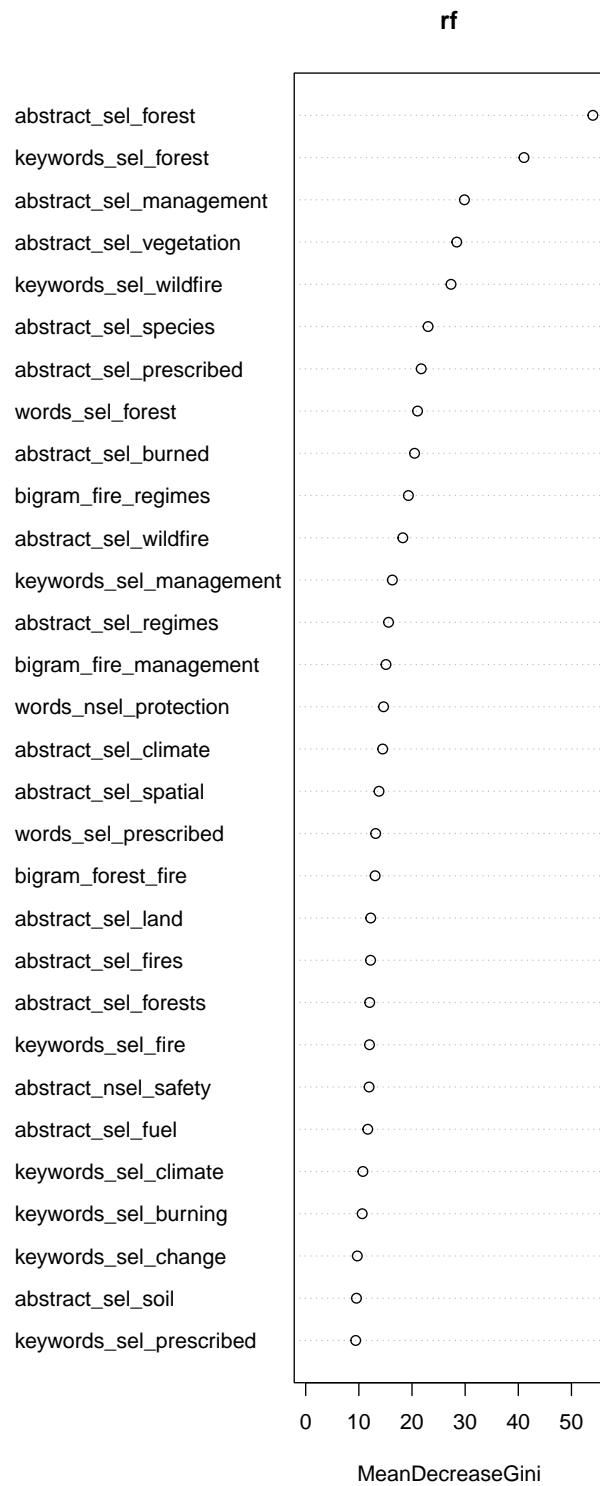
			0%
=====			11%



thresh	auc	recall	prec
0.439	0.921	0.927	0.922
0.379	0.924	0.951	0.903
0.527	0.913	0.917	0.914
0.423	0.931	0.945	0.924
0.456	0.926	0.941	0.916
0.484	0.924	0.930	0.923
0.417	0.930	0.943	0.921
0.510	0.929	0.912	0.949
0.452	0.930	0.919	0.945
0.450	0.932	0.931	0.937

thresh_avg	auc_avg	recall_avg	prec_avg	thresh_std	auc_std	recall_std	prec_std
0.454	0.926	0.931	0.925	0.044	0.006	0.013	0.014

1.5.3 Variable importance



List of the top-50 features by decreasing order of Mean Decrease in Gini Index:

MeanDecreaseGini	var_name
54.03	abstract_sel_forest
41.07	keywords_sel_forest
29.86	abstract_sel_management
28.42	abstract_sel_vegetation
27.34	keywords_sel_wildfire
23.02	abstract_sel_species
21.73	abstract_sel_prescribed
21.04	words_sel_forest
20.48	abstract_sel_burned
19.32	bigram_fire_regimes
18.26	abstract_sel_wildfire
16.30	keywords_sel_management
15.58	abstract_sel_regimes
15.10	bigram_fire_management
14.65	words_nsel_protection
14.48	abstract_sel_climate
13.79	abstract_sel_spatial
13.16	words_sel_prescribed
13.06	bigram_forest_fire
12.23	abstract_sel_land
12.20	abstract_sel_fires
12.04	abstract_sel_forests
12.00	keywords_sel_fire
11.93	abstract_nsel_safety
11.68	abstract_sel_fuel
10.76	keywords_sel_climate
10.63	keywords_sel_burning
9.71	keywords_sel_change
9.55	abstract_sel_soil
9.40	keywords_sel_prescribed
9.40	bigram_prescribed_fire
9.15	bigram_fire_regime
8.79	abstract_sel_landscape
8.22	abstract_nsel_thermal
8.22	abstract_sel_weather
7.78	keywords_nsel_protection
7.66	keywords_sel_ecosystem
7.41	abstract_nsel_protection
7.22	words_sel_weather

MeanDecreaseGini	var_name
7.09	bigram_fire_weather
6.57	abstract_sel_effects
6.53	abstract_sel_regime
6.37	bigram_forest_fires
6.21	abstract_sel_ecological
6.07	words_sel_regimes
6.02	abstract_nsel_heat
5.96	words_sel_management
5.75	bigram_fire_protection
5.60	abstract_sel_fire
5.23	abstract_sel_wildfires

1.5.4 Simplified model with best feature subset

Re-train the model but now with the top-50 best set of features based on the importance rank:

Call:

```
randomForest(x = train_tb %>% select(all_of(imp_vars)), y = train_tb %>% pull(included
```

```
  Type of random forest: classification
```

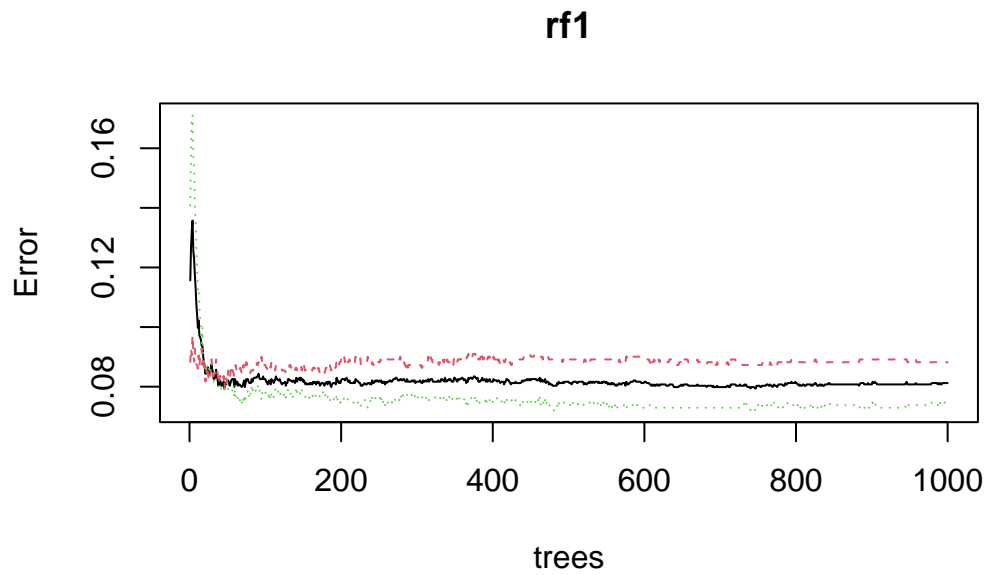
```
    Number of trees: 1000
```

```
No. of variables tried at each split: 12
```

```
      OOB estimate of  error rate: 8.12%
```

Confusion matrix:

```
      0      1 class.error
0 982    95 0.08820799
1  86 1066 0.07465278
```



1.6 Full model

1.6.1 Optimize cut-off value

thresh	auc	recall	prec
0.492	0.941	0.951	0.937

1.6.2 Predict class labels for the entire dataset

Predict class labels using the full model and also the optimized cut-off

Predicted class percentages:

```
pred_class
  0      1
37.2 62.8
```