

# WEB PAGE CLASSIFICATION

100042773 – Can Okan TAŞKIRAN

100042970 – Alperen KÖYLÜ

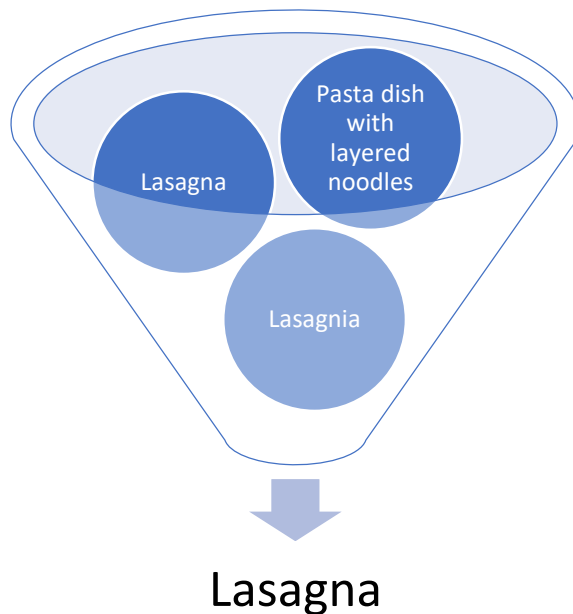
# 1. Introduction

# What is the problem ?

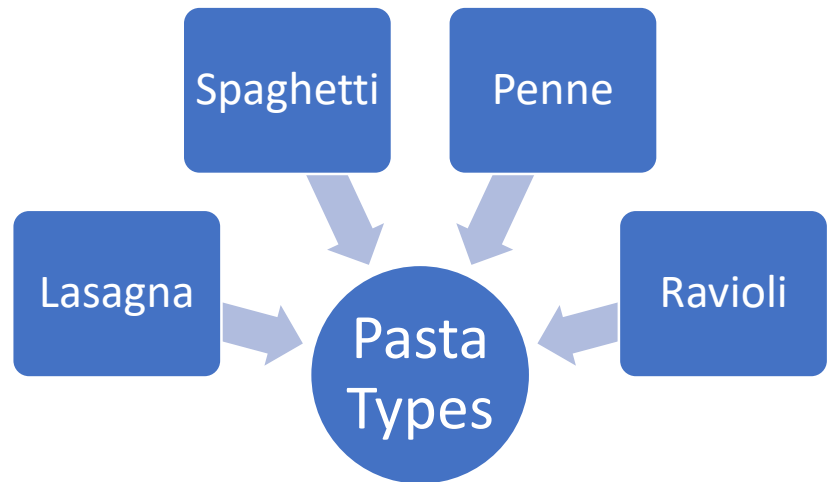
- Past decade we have witnessed an explosive growth on the Internet, with millions of web pages on millions of topics.
- But the Internet is not a suitable tool for locating or organizing the mass of information.
- Tools like search engines assist users in locating information on the Internet.
- They perform excellently in locating but provide limited ability in organizing the web pages.

# Difference between locating and classification

What **Google** do ...



What we are trying ...



# What is web page classification ?

- Web page classification, aka categorization, defined as the task of determining whether a web page belongs to a category or categories.
- Let  $\mathbf{C} = \{\mathbf{c1}, \dots, \mathbf{cK}\}$  be a set of predefined categories,  $\mathbf{D} = \{\mathbf{d1}, \dots, \mathbf{dN}\}$  be a set of web pages to be classified, and  $\mathbf{A} = \mathbf{D} \times \mathbf{C}$  be a decision matrix.

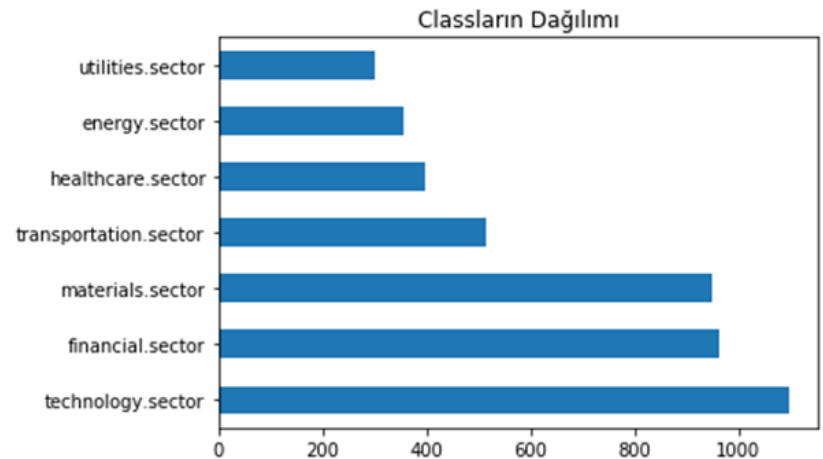
# A diagram for showing the decision matrix

Web Pages	Categories				
	C1	...	Cj	...	Ck
D1	A11	...	A1j	...	A1k
⋮	...	...	...	...	...
Di	Ai1	...	Aij	...	Aik
⋮	...	...	...	...	...
Dn	An1	...	Anj	...	Ank

## 2. Experimental Setup

# Our working space

- 7 classes
  - Basic materials sector
  - Energy sector
  - Financial sector
  - Healthcare sector
  - Technology sector
  - Transportation sector
  - Utilities sector
- 4581 individual html files
- More than 1461640 dimensions in matrix (no reduction or selection algorithm applied on)





# Pure HTML view

```
<!DOCTYPE HTML PUBLIC "-//IETF//DTD HTML 2.0//EN">
<HTML>
<HEAD>
<TITLE>Kaydon Bearings</TITLE>
<META name="description" content="Kaydon bearings are the designer's choice for difficult and sensitive applications in material handling, construction equipment, robotic and other key industrial positions.">
<META name="keywords" content="Bearings, Reali-Slim, Thin Section, Ball, Roller, Custom, Assembly, Hybrid, Stainless-steel, Vacuum, Robots, Large bore, Turntable, Slewing Rings, Worm Drive, Construction Equipment, Manlifts">

</HEAD><body background="/kaydon/backgrounds/bgmenu.gif">
<table width="500">
  <td width="160" align="left" valign="top">
<center><A HREF="/kaydon/default.htm"><IMG SRC="/kaydon/graphics/logo100.gif" border="none"></A>
<IMG SRC="/kaydon/bearings/graphics/bearings3.jpg" width="100">
<A HREF="/kaydon/bearings/default.htm"><IMG SRC="/kaydon/bearings/graphics/bearings.gif" border="none"></A></center>
</td>
<td align="center"><A HREF="/kaydon/bearings/index.htm">Bearings Site Index</A></center></td>
</tr>
<tr>
<td align="center"><IMG SRC="/kaydon/bearings/graphics/bearingslogo.gif"></center>
<table width="100%" border="1">
  <tr>
<td align="center"><IMG SRC="/kaydon/bearings/graphics/bearingspix.jpeg"></center>
</td>
</tr>
<tr>
<td align="center"><B>2860 McCracken St<BR>Muskegon, MI 49441<BR>(616) 755-3741<BR>FAX: (616) 759-4102<BR>Telex: 228436<BR><BR><hr>Engineer Services <BR>Hotline <BR>(800) 514-3066</td>
</tr>
</table>
</tr>
<tr>
<td align="center">
<P>
<A HREF="/kaydon/bearings/cap.htm">Capabilities</A> |
<A HREF="/kaydon/bearings/app.htm">Applications</A> |
<A HREF="/kaydon/bearings/products.htm">Products</A> |
<A HREF="/kaydon/bearings/interface.htm">Interface</A> |
<A HREF="/kaydon/bearings/faqs.htm">FAQs</A></P>
</td>
</tr>
<tr>
<td align="center">
<B>Kaydon bearings</B> are the designer's choice for difficult and sensitive applications in material handling, construction equipment, robotic and other key industrial positions. Carefully engineered solutions are available promptly and completely. Prototype
<center>
Other Kaydon Divisions at your service:<BR>
<A HREF="/kaydon/cooper/default.htm">Cooper Split Roller Bearings</A> |
<A HREF="/kaydon/electro/default.htm">Electro-Tec</A> |
<A HREF="/kaydon/filtration/default.htm">Filtration</A> |
<A HREF="/kaydon/fluid/default.htm">Fluid Power</A> |
<A HREF="/kaydon/iti/default.htm">Industrial Tectonics Inc</A> |
<A HREF="/kaydon/rings&seal/default.htm">Kaydon Rings & Seal</A>
```

# How to interpret HTML files ?

- External wrote C# code that removes html tags and styles in pages, and gives you only texts.
- BeautifulSoup4
  - P tags
  - Meta tags
  - All tags
- P.S. : By the way cross-validation applied all algorithms.

# Challenges

## Sparse Matrix

Index	bearing	company	contract	diving	foil	home	information	jrm	management	oil	paper	power	service	software	system	technology	treatment	united	wastewater	water	class
4417	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4422	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4454	0	0.8663834	0	0.8...	0	0	0	0	0	0.294874	0	0.8713...	0.37727	0	0.256557	0.8718893	0	0	0	0	1
4457	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4479	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4481	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4498	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4514	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4516	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4517	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1
4545	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4576	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
25	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
37	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
48	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
42	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
43	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

## Only P Tags

Index	bearing	company	contract	diving	foil	home	information	jrm	management	oil	paper	power	service	software	system	technology	treatment	united	wastewater	water	class
4417	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4422	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4454	0	0.8663834	0	0.8...	0	0	0	0	0	0.294874	0	0.8713...	0.37727	0	0.256557	0.8718893	0	0	0	0	1
4457	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4479	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4481	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4498	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4514	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4516	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4517	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	1
4545	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
4576	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
7	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
21	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
23	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
25	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
37	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
48	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
42	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
43	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

# Solution

- Final matrix that all learning algorithm applied on

index	company	content	date	genre	html	image	inc	information	last	length	modified	new	mov	product	server	service	system	text	type	year	class
619	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
884	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
1020	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
1254	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
1353	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
1396	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
1709	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
1942	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
2015	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
3248	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
3956	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
4121	0	0.33371	0.337	0.338	0.333920	0	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3338	0.33	0	0
767	0	0.271125	0.411	0.275	0.135651	0	0	0.704983	0.150011	0.157488	0.100515	0	0.187	0	0.1358	0	0	0.1356	0.13	0	0
1714	0	0.0038	0.0039	0.003	0.00385031	0.0143	0	0	0	0.0447964	0	0.99	0.005	0	0.0038	0	0	0.0154	0.00	0	1
968	0.905834	0.0099	0.005	0.010	0.00409202	0	0.411	0.00064786	0.00561	0.00579562	0.00590701	0.03	0.006	0.0187	0.0050	0.0012549	0	0.0048	0.00	0.0374205	1
297	0.953986	0.0127	0.012	0.012	0.0127669	0	0.025	0	0	0	0	0.07	0.017	0.2878	0.0127	0	0	0.0127	0.01	0	1
1205	0	0.0259	0.013	0.026	0.998176	0	0	0	0.01510	0.0150501	0.0153394	0	0.017	0.0245	0.0129	0	0	0.0129	0.01	0	1
4527	0.45793	0.0264	0.014	0.028	0.0142244	0	0	0.0492629	0.04878	0.0165142	0.0168316	0.02	0.019	0	0.0142	0.0257256	0	0.0142	0.01	0.746501	1
2576	0.647973	0.0525	0.017	0.035	0.0173446	0.0644	0	0.0000034	0.101018	0.0201367	0.0205237	0.05	0.023	0.0325	0.0173	0.0313686	0	0.0173	0.03	0.736009	1
4326	0.0636208	0.0178	0.018	0.036	0.0178812	0	0	0.0019527	0	0	0	0.16	0.024	0.4705	0.0178	0.380009	0.702	0.0178	0.01	0	1

# Lemmitazing vs Stemming

## Lemmitazing

index	company	content	date	gmt	html	image	pic	information	last	length	modified	new	nov	product	server	service	system	test	type	year	class
818	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
884	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
1808	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
1214	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
1355	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
8396	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
1709	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
1942	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
2815	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
3248	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
3856	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
4121	0	0.33371	0.337	0.338	0.333928	0	0	0	0	0	0	0.461	0	0.3345	0	0	0.3339	0.33	0	0	0
767	0	0.271125	0.411	0.275	0.135851	0	0	0.704881	0.158811	0.157488	0.160515	0	0.187	0	0.1376	0	0.1356	0.13	0	0	0
3714	0	0.0838	0.0839	0.083	0.08383851	0.0143	0	0	0	0.0447804	0	0.98	0	0.005	0	0.0038	0	0.0154	0.08	0	1
968	0.968834	0.0099	0.005	0.018	0.00499282	0	0.41	0.00884768	0.00581	0.00579582	0.00580781	0.01	0.006	0.0187	0.0058	0.0012549	0	0.0045	0.08	0.0374261	1
287	0.953088	0.0127	0.012	0.012	0.0127669	0	0.025	0	0	0	0	0.07	0.017	0.0878	0.0127	0	0	0.0127	0.01	0	1
12495	0	0.0288	0.013	0.026	0.008139	0	0	0.01318	0.0188081	0.0133394	0	0.017	0.0243	0.0128	0	0	0.0128	0.01	0	1	0
45117	0.65789	0.0184	0.014	0.028	0.0142144	0	0	0.0492829	0.04978	0.0105182	0.0108318	0.01	0.018	0	0.0182	0.01257256	0	0.0182	0.01	0.7465881	1
2576	0.647971	0.0128	0.017	0.035	0.0173446	0.0044	0	0.0000934	0.181818	0.0201167	0.0205237	0.09	0.023	0.0325	0.0173	0.0013688	0	0.0173	0.01	0.7368869	1
4326	0.0036288	0.0178	0.018	0.036	0.0178812	0	0	0.0613927	0	0	0	0.16	0.024	0.4793	0.0278	0.388889	0.767	0.0278	0.01	0	1

## Stemming

index	company	content	corpor	date	gmt	html	imag	pic	inform	last	length	manag	modifi	new	nov	oper	product	server	servic	system	test	thi	type	use	year	class
81	0	2	4	1	2	1	0	0	1	2	2	0	1	0	1	1	0	1	7	1	1	1	1	0	0	0
92	1	2	0	1	1	1	0	0	1	1	1	0	1	0	0	0	1	0	0	2	1	1	0	0	0	0
180	2	1	0	1	1	1	0	0	1	0	0	0	0	3	0	0	0	1	4	0	1	0	1	0	0	0
114	0	2	0	1	2	1	0	0	1	1	0	0	1	0	0	0	1	0	0	1	0	1	0	0	0	0
119	3	2	0	1	2	1	0	0	1	1	1	0	1	0	0	2	3	1	2	1	1	1	1	2	0	0
132	1	2	0	1	2	1	0	0	0	1	1	0	1	0	0	1	0	1	1	0	2	0	1	0	0	0
139	7	2	1	1	1	1	0	0	1	1	0	1	1	0	0	0	1	2	5	1	2	1	0	0	0	0
181	1	2	0	2	2	1	0	0	1	1	0	1	1	1	0	0	1	0	0	1	2	1	0	0	0	0
198	0	1	0	1	1	1	0	0	2	0	0	0	0	0	0	0	2	0	0	1	0	1	0	1	0	0
211	0	2	0	1	2	1	0	0	0	1	1	0	1	0	0	0	1	0	0	1	1	1	0	0	0	0
214	0	1	0	1	1	1	0	0	0	0	0	0	0	1	0	0	0	1	4	1	1	0	1	0	0	0
219	0	2	0	1	2	1	0	0	0	2	1	0	1	0	0	0	0	1	2	0	1	0	1	0	0	0
221	0	3	0	0	0	1	0	0	0	1	0	0	0	0	0	0	2	0	0	1	1	1	0	0	0	0
225	0	2	4	1	2	1	0	0	1	1	1	0	1	0	0	0	1	2	0	1	1	1	0	0	0	0
235	1	3	2	1	2	1	0	0	7	0	1	1	1	2	0	0	1	4	0	1	2	1	0	0	0	0
242	0	2	0	1	2	1	0	0	1	1	1	0	1	0	0	0	1	1	0	1	0	1	0	1	0	0
282	6	2	4	1	3	1	0	0	2	0	1	1	0	1	0	2	4	1	3	0	1	2	1	0	0	0
123	15	2	7	1	2	1	0	0	20	2	1	1	1	1	0	0	2	1	11	0	1	2	1	4	0	0
371	0	2	0	1	3	1	0	0	1	0	1	1	0	0	0	0	1	0	1	2	0	1	0	0	0	0
383	0	2	0	1	2	1	0	0	1	1	1	0	1	0	0	0	2	0	0	2	1	1	0	0	0	0

P.S. : Stopwords and regex are eliminated before we process on data on every algorithm.

# Future representation

- TF-IDF
- Count Vectorizer + TF-IDF
- N Gram
  - BiGram

# TOP 5 COUNT+TF-IDF

MAX_FEATURES	MIN_DF	MAX_DF	DTREE
100	600	250	0.713787086
90	600	90	0.712041885
100	600	310	0.710296684
90	500	230	0.709424084
90	500	310	0.709424084

MAX_FEATURES	MIN_DF	MAX_DF	KNN
70	750	90	0.620418848
80	750	90	0.620418848
90	750	90	0.620418848
100	750	90	0.620418848
60	750	90	0.620418848

MAX_FEATURES	MIN_DF	MAX_DF	BAYES
90	450	110	0.533158813
90	450	510	0.533158813
90	450	150	0.533158813
90	450	310	0.533158813
90	450	330	0.533158813

Obtained from **accuracy\_VECTORIZER&TFIDF.csv** file and this file has been created after 23 hours of live working code.

# TOP 5 TF-IDF

MAX_FEATURES	MIN_DF	MAX_DF	DTREE
140	450	70	0.725130890
130	450	90	0.721640489
140	450	110	0.721640489
140	450	250	0.721640489
140	450	370	0.719895288

MAX_FEATURES	MIN_DF	MAX_DF	KNN
150	450	90	0.623909250
140	450	90	0.623909250
70	750	90	0.620418848
100	750	90	0.620418848
90	750	90	0.620418848

MAX_FEATURES	MIN_DF	MAX_DF	BAYES
130	450	270	0.538394415
130	450	110	0.538394415
130	450	470	0.538394415
130	450	390	0.538394415
130	450	370	0.538394415

Obtained from **accuracy\_TFIDFonly.csv** file and this file has been created after 11 hours of live working code.



# Features Selection and Feature Extraction

- Feature Selection

- Filter Methods

- Correlation
    - Information Gain (Mutual Inf.)
    - Relief

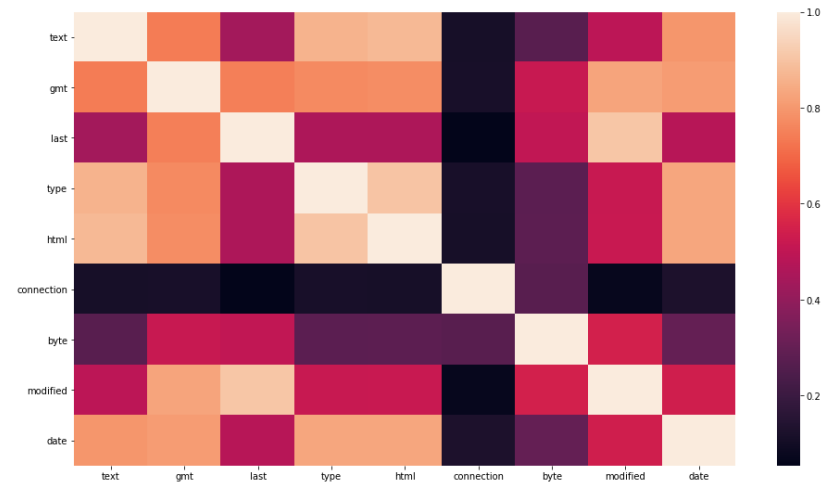
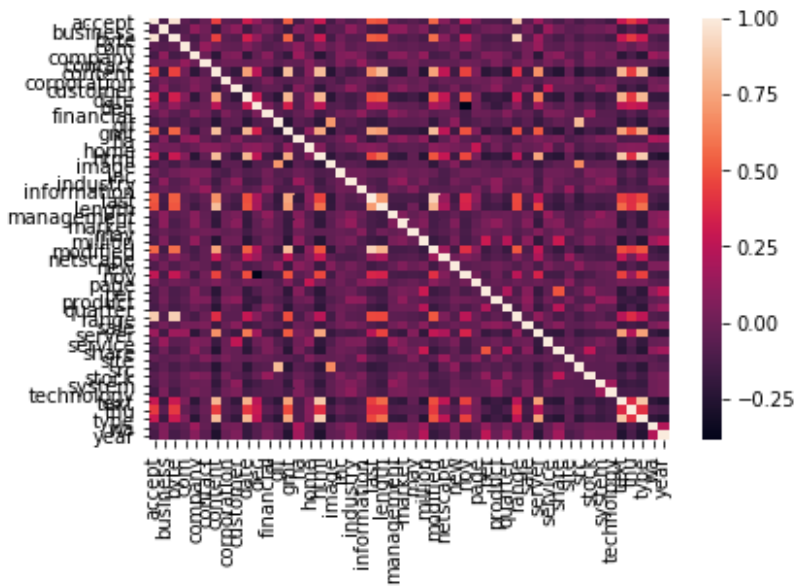
- Wrapper Methods

- Sequential Feature Selection

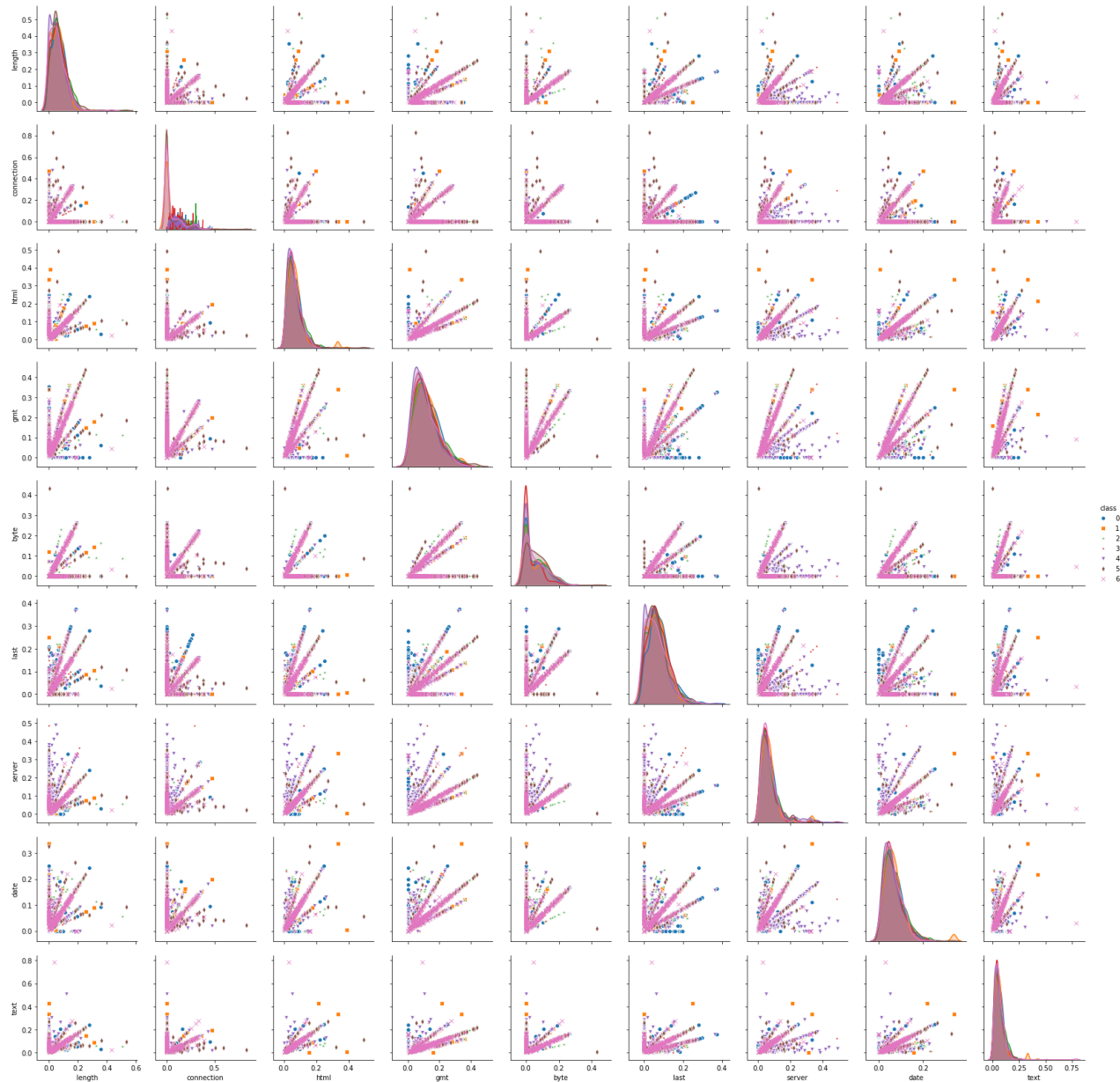
- Feature Reduction

- Principal Component Analysis (PCA)
  - Linear Discriminant Analysis (LDA)
  - Latent Semantic Analysis (LSA)

# Correlation matrix



# Cross of Features and Correlation



# Latent Semantic Analysis (SVD)

- Unsupervised
- Decompose it into 3 separated matrix.
- Term-Document Matrix

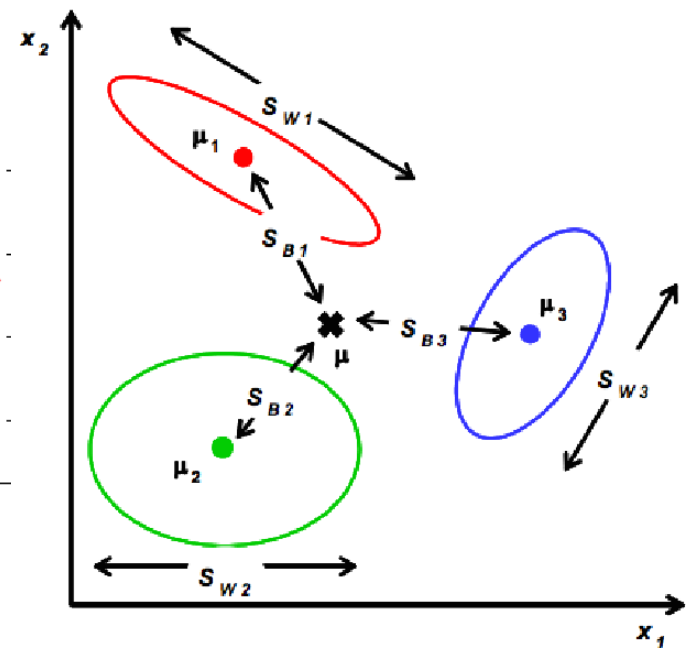
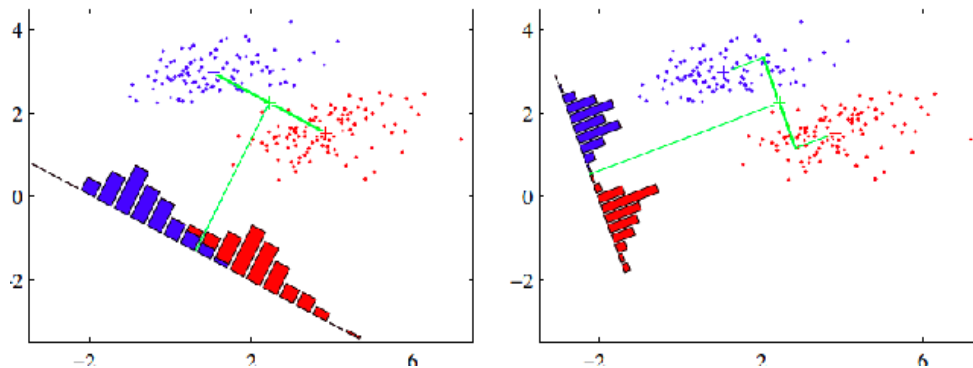
$$A = U S V^T$$

$$A_{m \times n} = U_{m \times m} S_{m \times n} V_{n \times n}^T$$

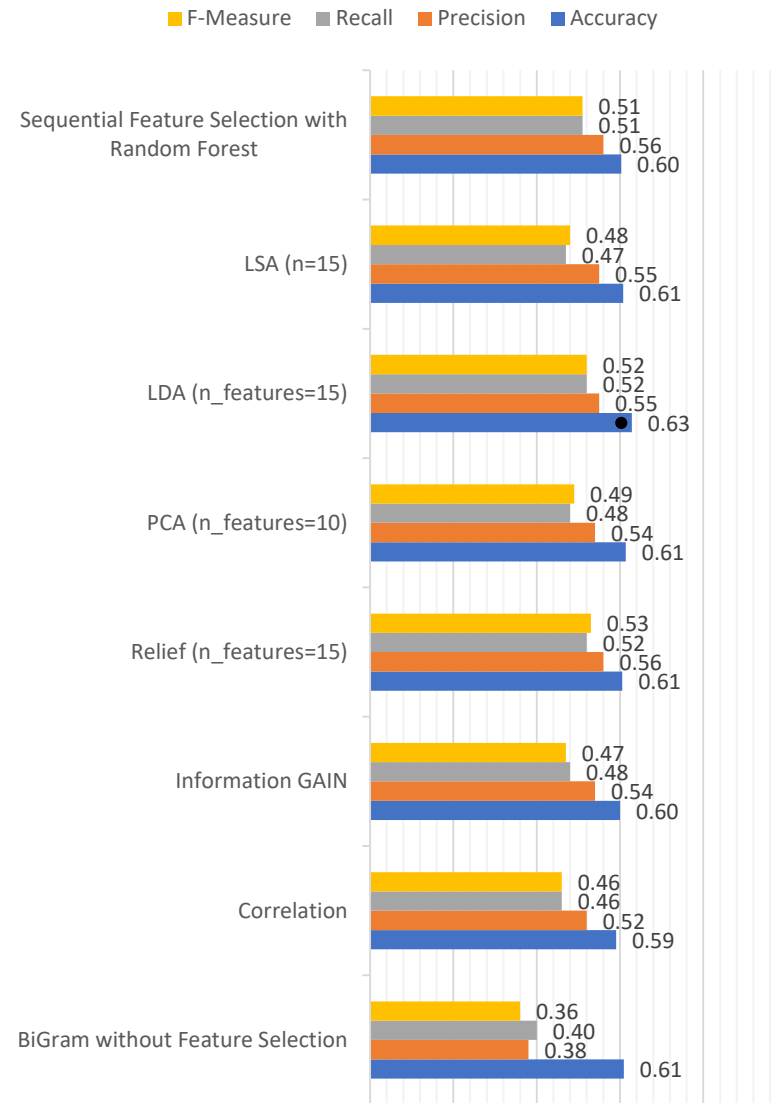
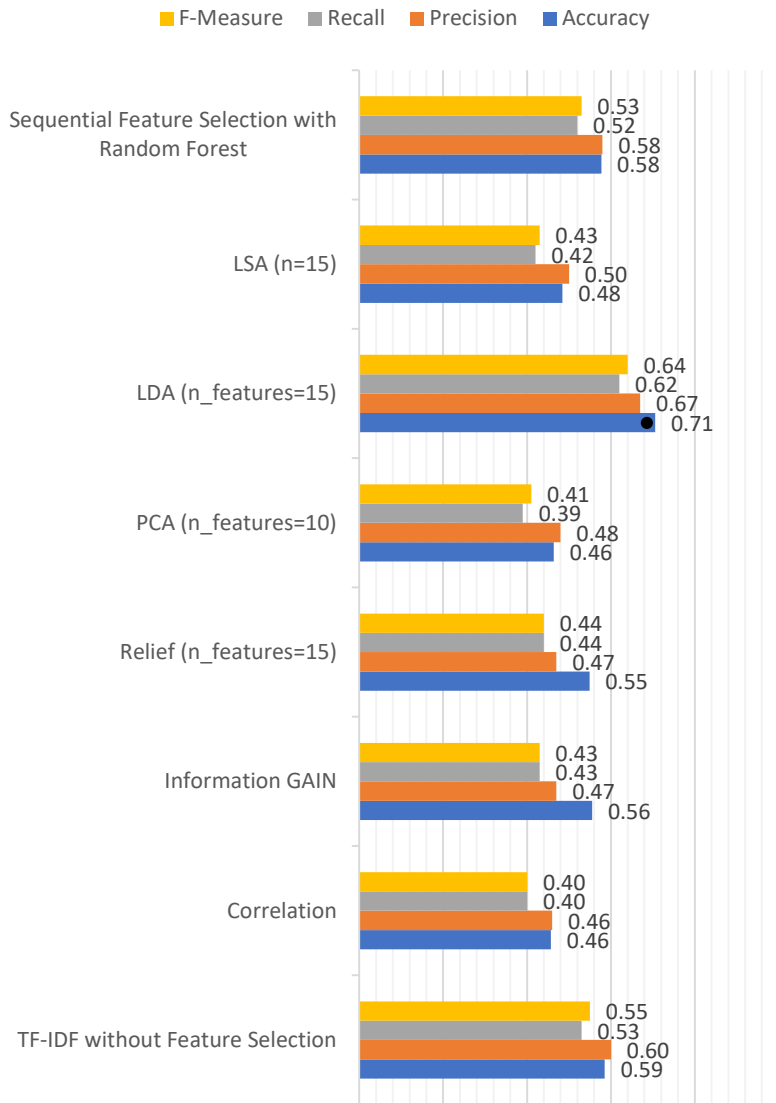
$$\begin{array}{c} \begin{array}{c} A \\ \left( \begin{array}{ccc} x_{11} & x_{12} & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{m1} & & x_{mn} \end{array} \right) \\ m \times n \end{array} = \begin{array}{c} \begin{array}{c} U \\ \left( \begin{array}{ccc} u_{11} & & u_{m1} \\ & \ddots & \\ u_{1m} & & u_{mm} \end{array} \right) \\ m \times m \end{array} \begin{array}{c} \begin{array}{c} S \\ \left( \begin{array}{ccc} \sigma_1 & & 0 \\ & \ddots & \\ 0 & & \sigma_r & \ddots & 0 \end{array} \right) \\ m \times n \end{array} \begin{array}{c} \begin{array}{c} V^T \\ \left( \begin{array}{ccc} v_{11} & & v_{1n} \\ & \ddots & \\ v_{n1} & & v_{nn} \end{array} \right) \\ n \times n \end{array} \end{array} \end{array}$$

# Linear Discriminant Analysis (LDA)

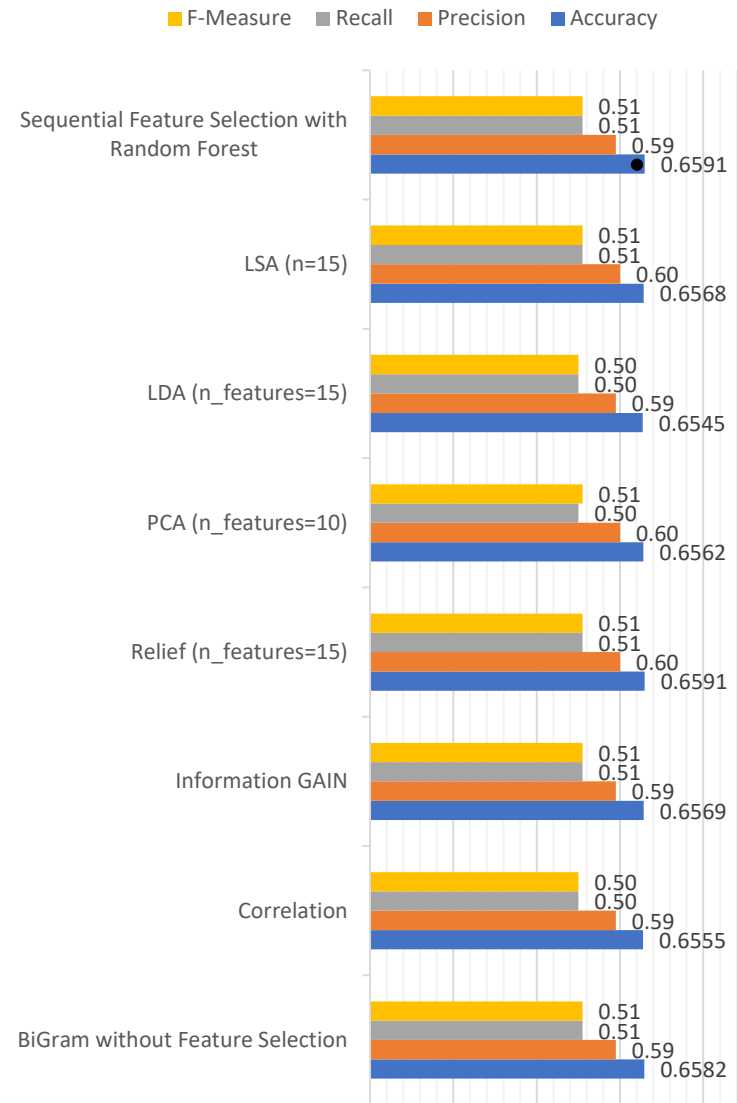
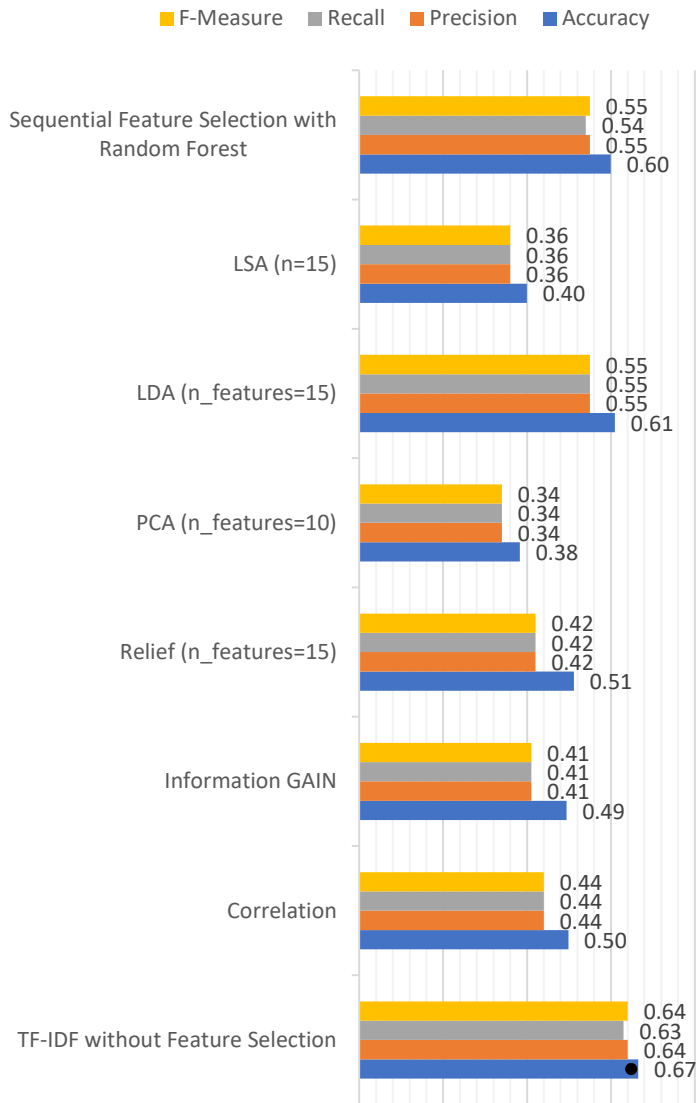
- Supervised
- Preserving class discrimination Information
- Maximize distance



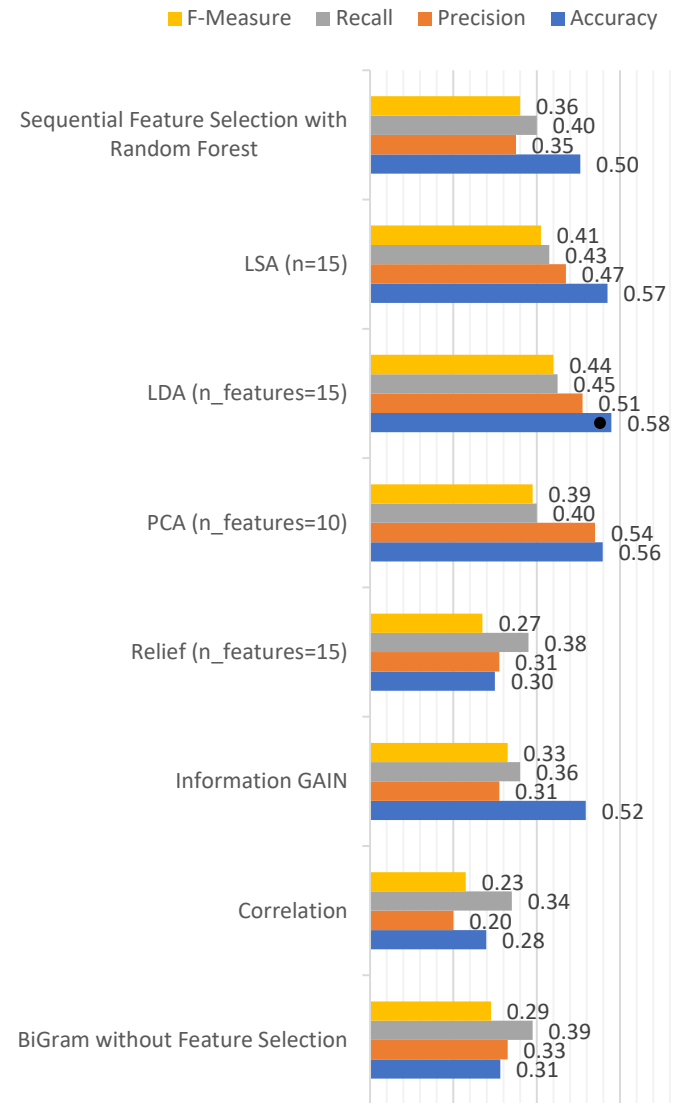
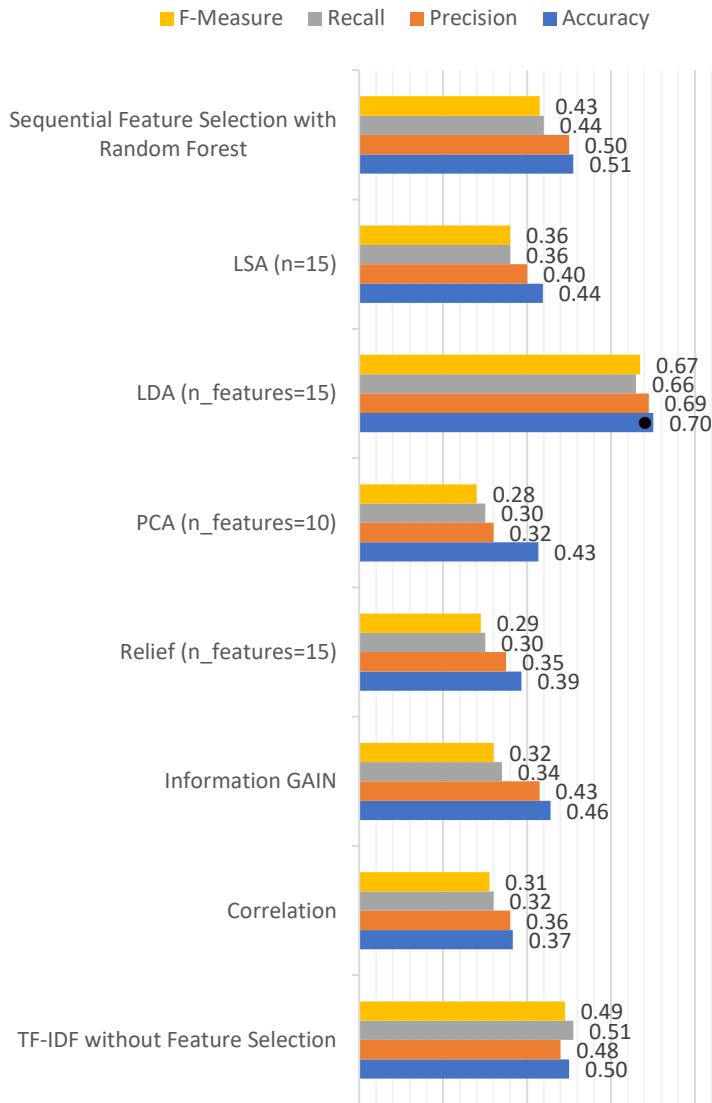
# Count Vectorizer vs. Bi Gram on KNN



# Count Vectorizer vs. Bi Gram on Decision Tree

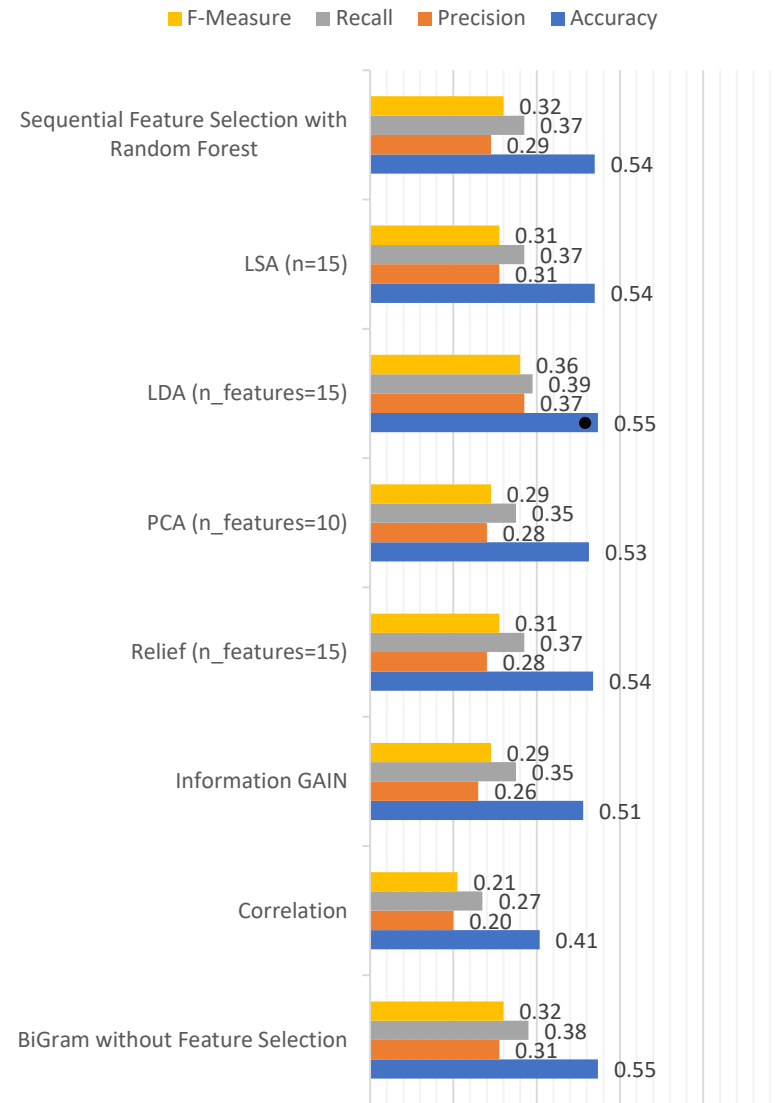
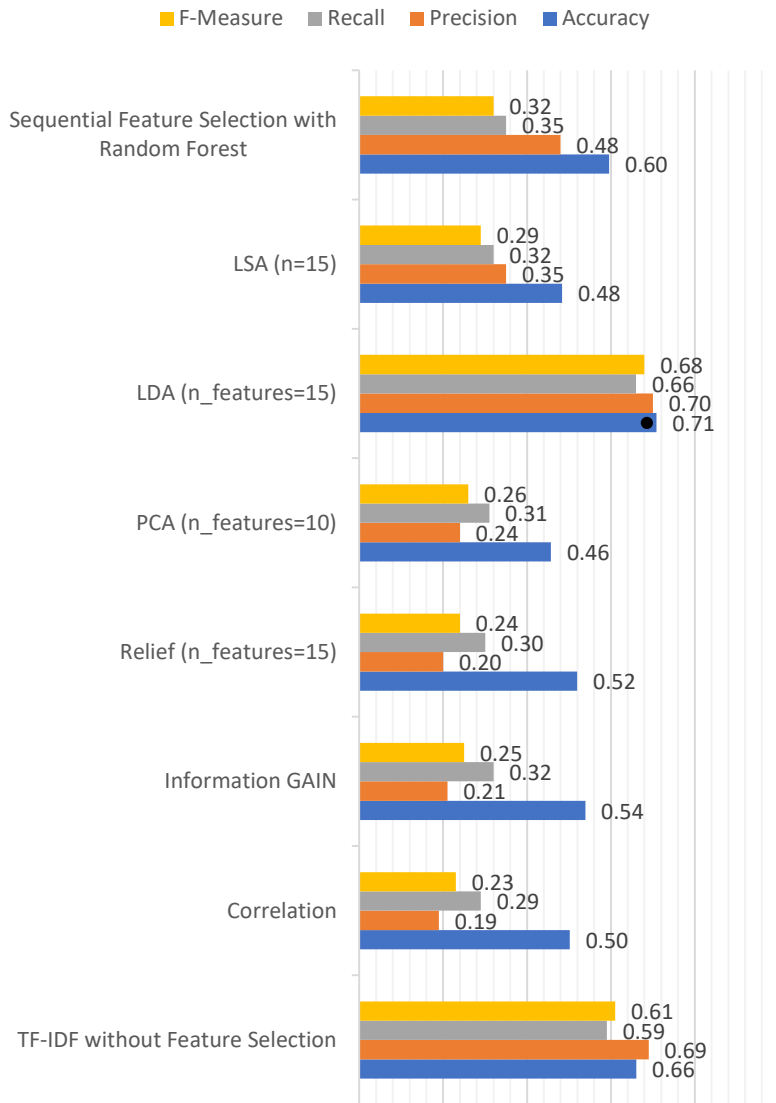


# Count Vectorizer vs. Bi Gram on Naive Bayes

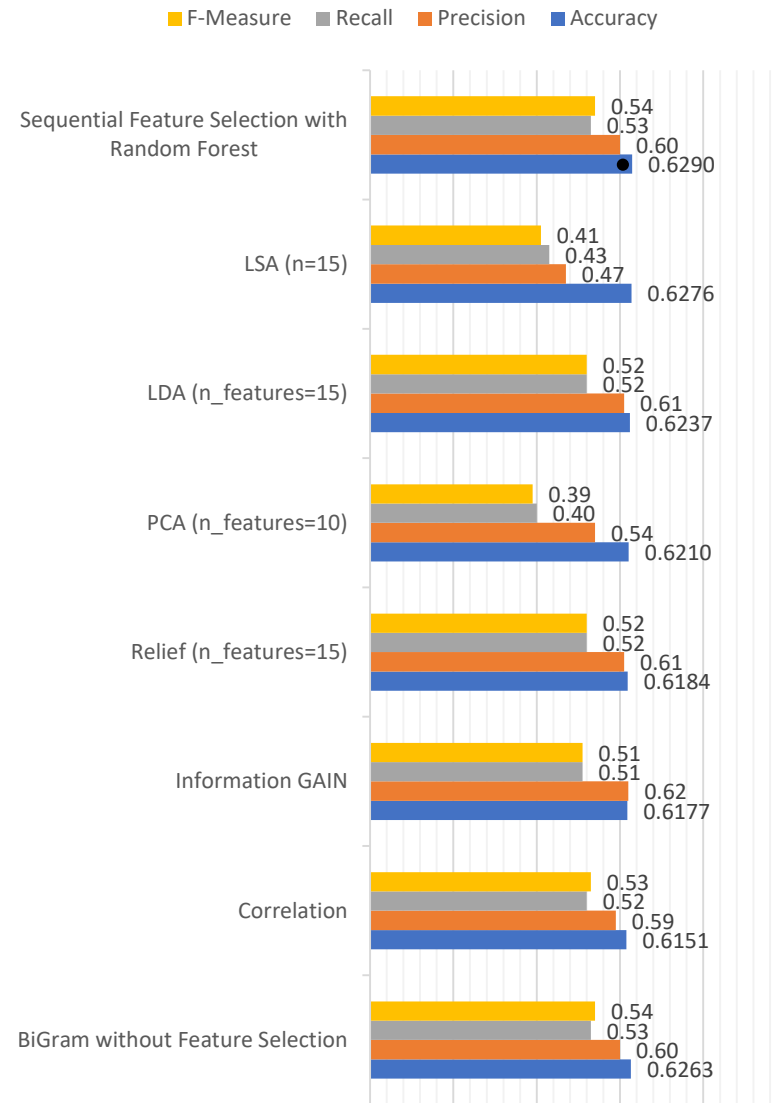
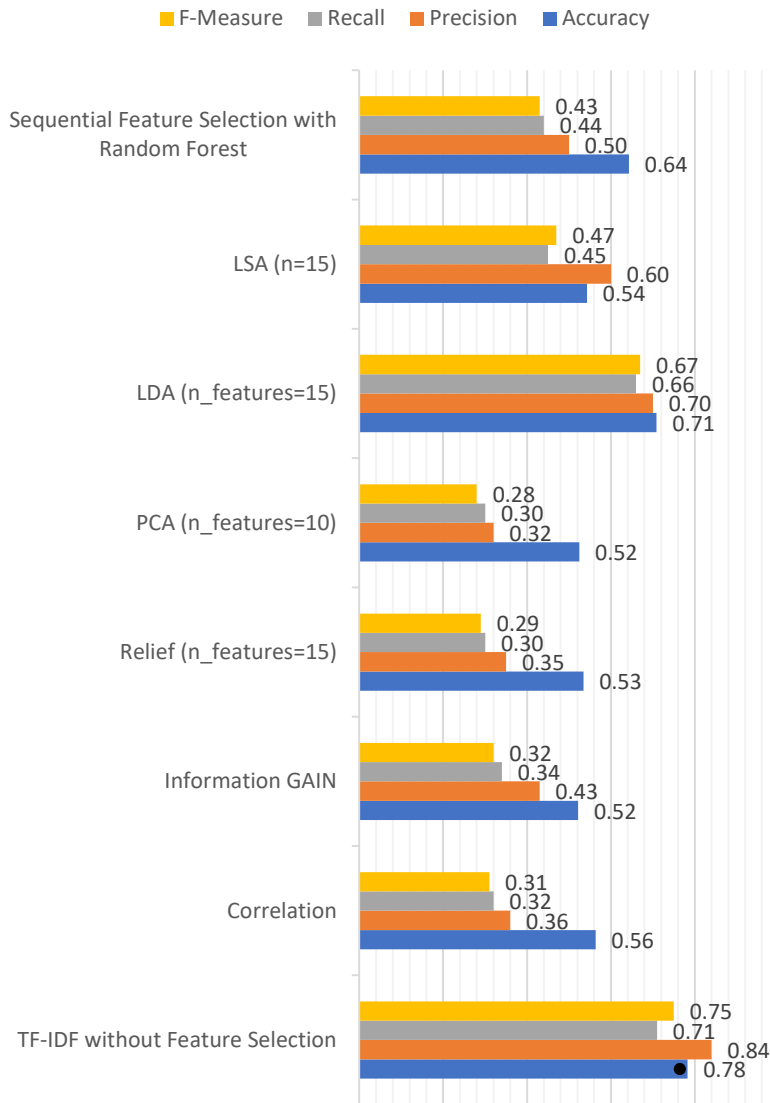




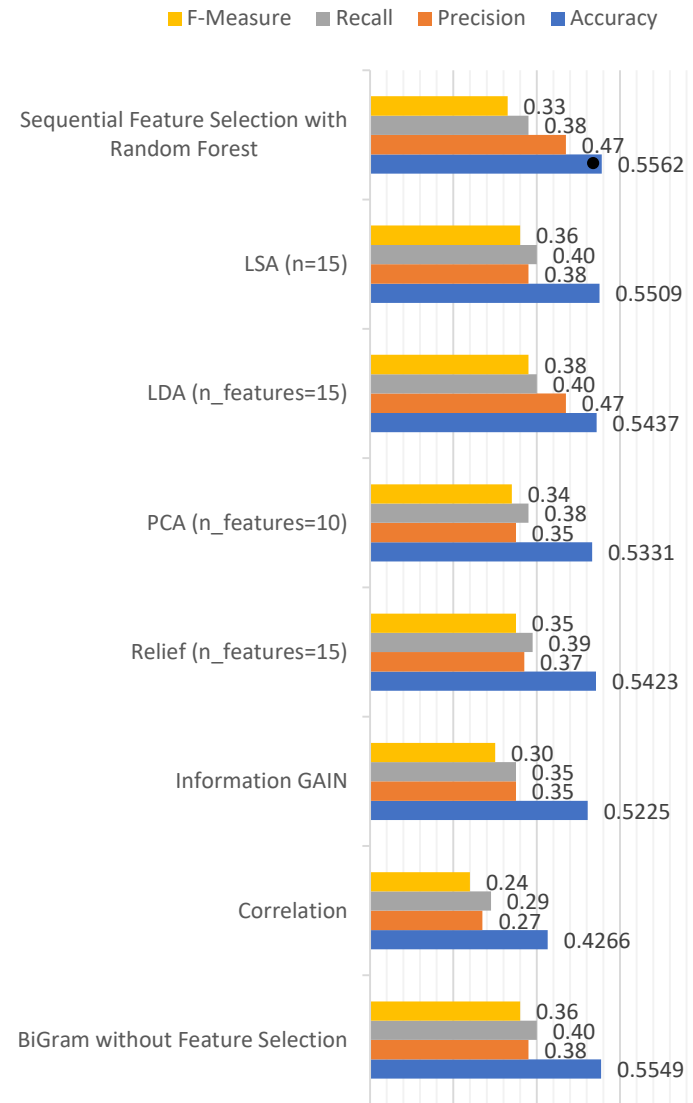
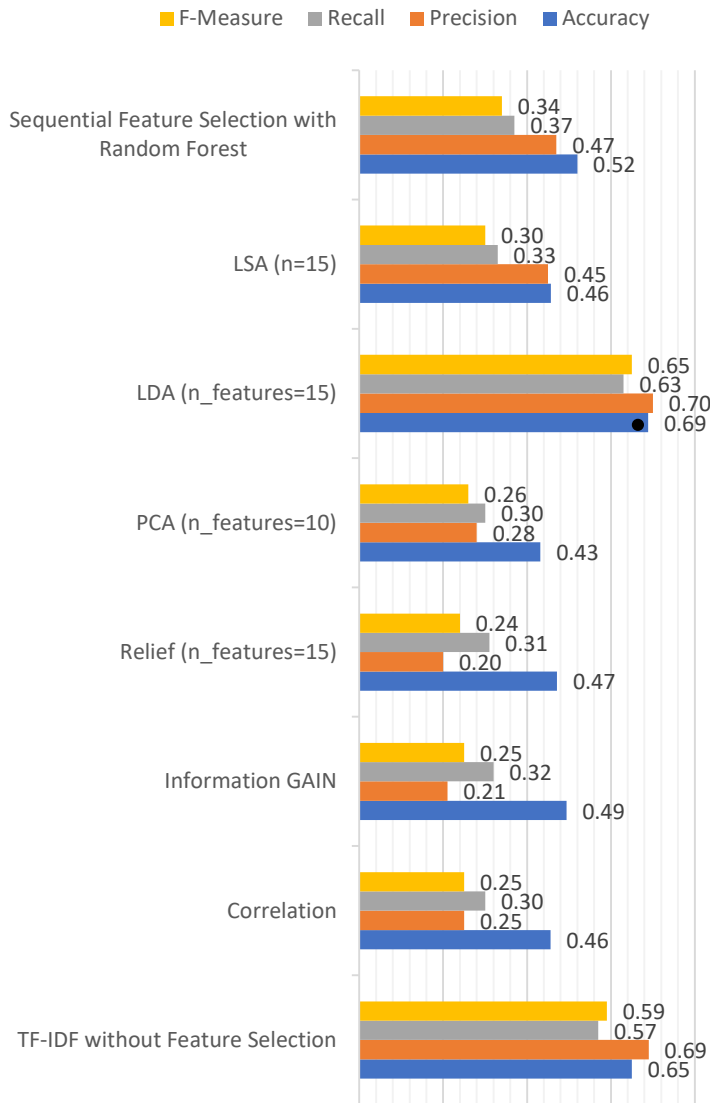
# Count Vectorizer vs. Bi Gram on SVM



# Count Vectorizer vs. Bi Gram on Random Forest

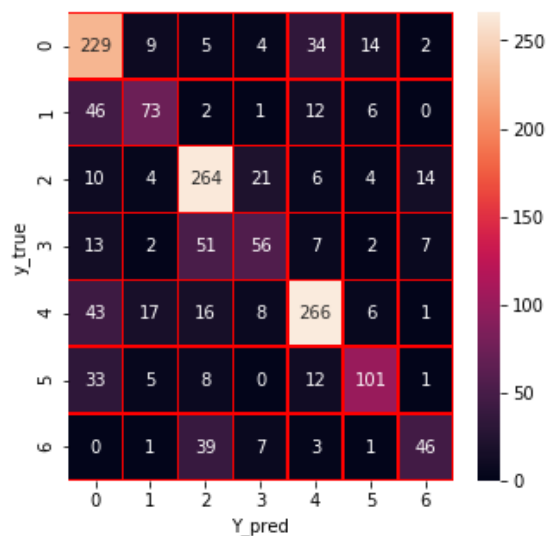


# Count Vectorizer vs. Bi Gram on Logistic Regression

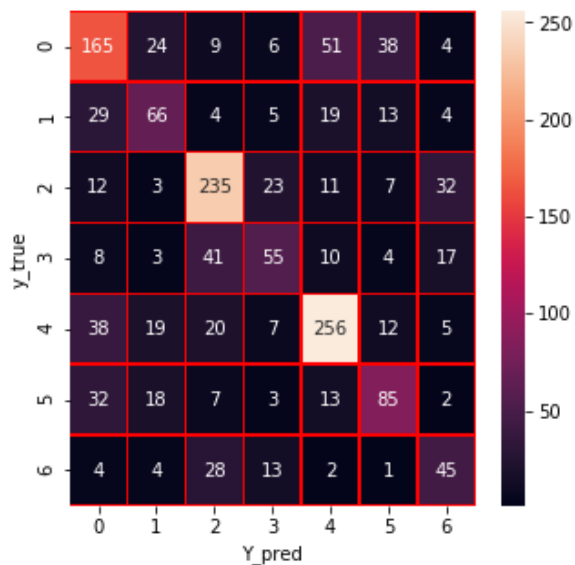


# Confusion Matrix (LDA)

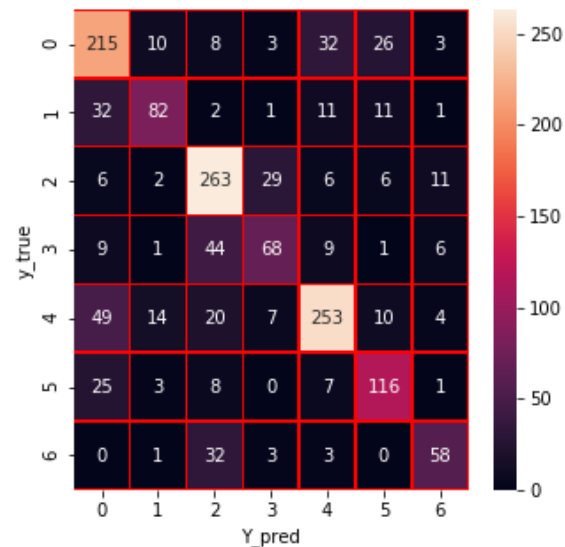
KNN



Decision Tree

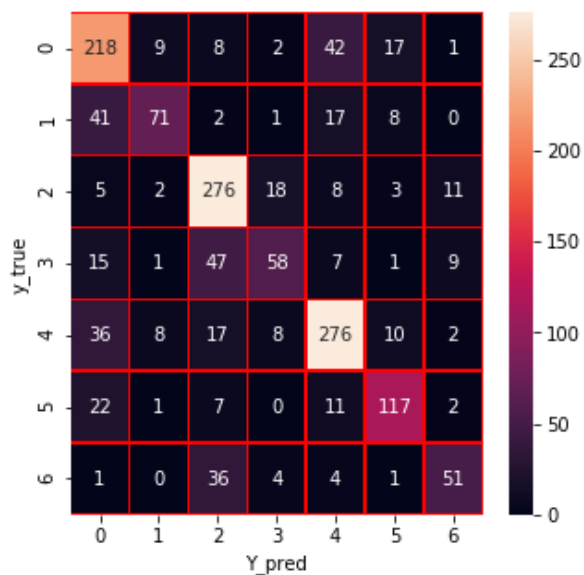


Naive Bayes

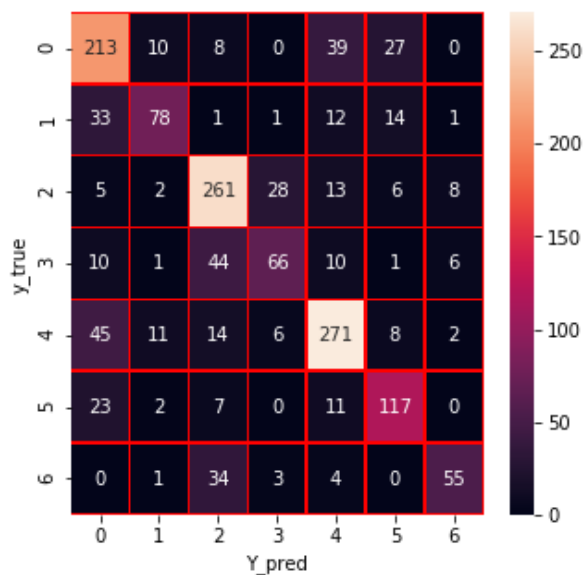


# Confusion Matrix (LDA)

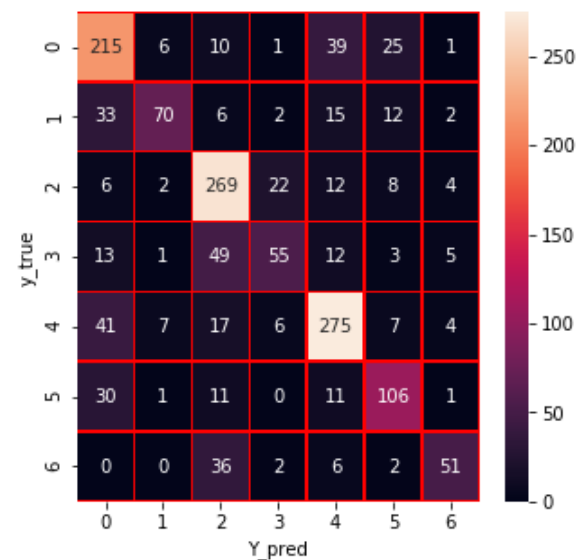
Random Forest



SVM



Logistic Regression



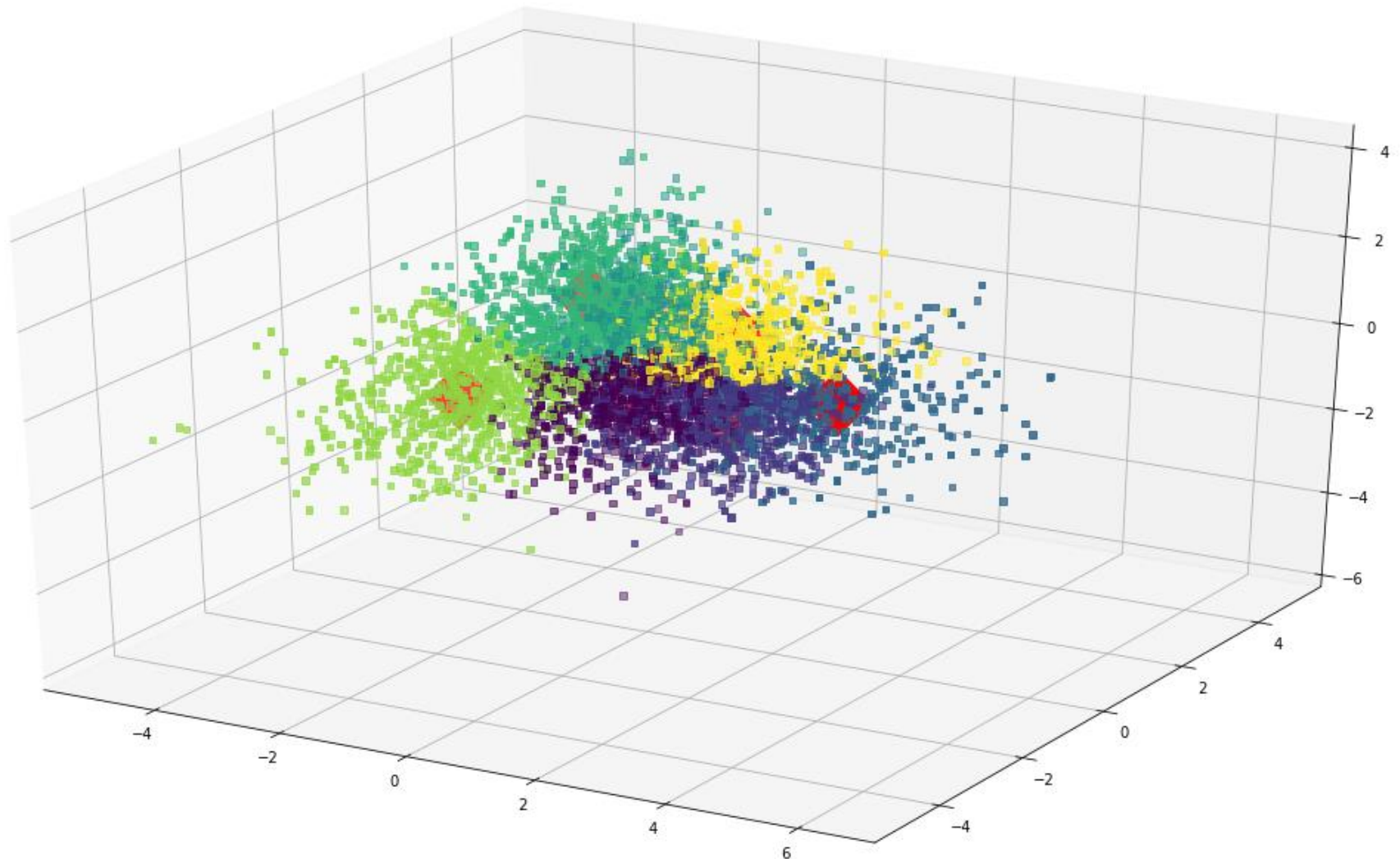
# Summary

- As we seen on the charts and matrixes, the most successful feature selection algorithm is TF-IDF without future selection in general.
- But, working with the large dimension brings us different kinds of problems.
  - As the dimensions increase, time complexity's impact will be huge on learning models (Kind of curse of dimensionality).
  - Irrevelant and outliers data.
  - Visualization.

# Summary

- So that; we decided to continue with the LDA algorithm.
- While we working on LDA, our features number decreased to six (thanks to data mining gods).
- We got nearly same results with the features of six instead of thousands. This fact shows us the power of LDA algorithm very strongly.

# Extra( K-Means, Hierarchical C.)



n\_clusters=7 - n\_components=3



# Dendrogram

