

CONTENT

- 01 INTRODUCTION
- 02 DATA
- 03 PROCESS
- 04 RESULTS ATTAINED + IMPACT
- 05 BUMPS + OBSTACLES
- 06 SUMMARY + QUESTIONS

WHY IS THIS IMPORTANT?



- **Main idea:** Investigating the integration of machine learning models in GitHub pull requests.
- **Process:** Analyzing a set of successful GitHub pull requests and identifying patterns and factors that contribute to their acceptance.
- **Final Goal:** Hope to improve the efficiency of software development and its processes.

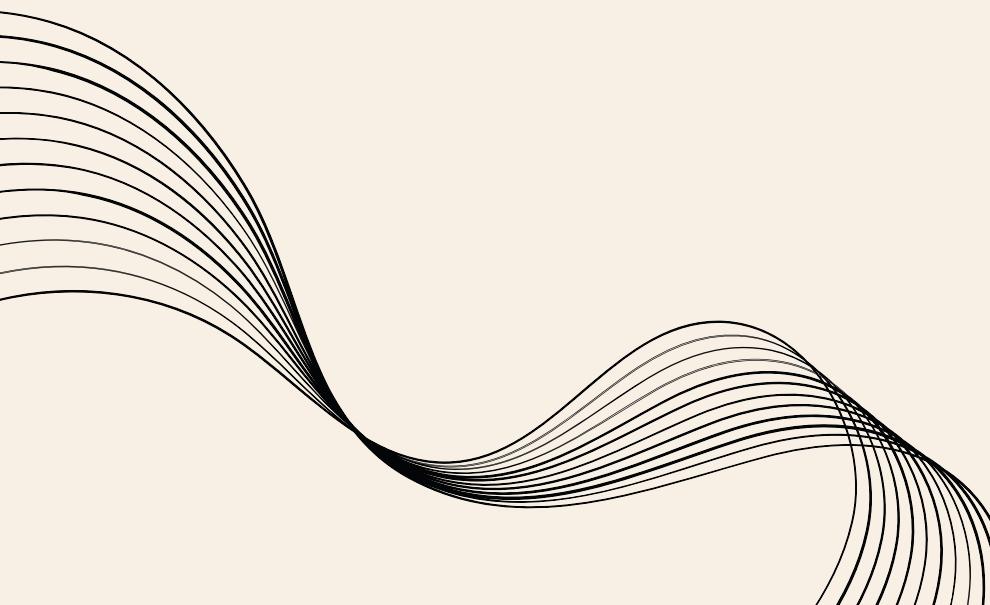
AN EMPIRICAL ANALYSIS OF PULL REQUESTS:

**EXPLORING THE INTEGRATION OF MACHINE LEARNING MODELS
IN SOFTWARE DEVELOPMENT**

**CLAUDIA FARKAS
PROFESSOR MALEKNAZ NAYEBI**

THE RESEARCH QUESTION

**How Can We Identify and Classify Messages
Contained in a Pull Request for Machine
Learning Initiatives?**



SOME BACKGROUND

https://github.com/tensorflow/tensorflow/pull/41178#discussion_r453906256

A screenshot of a GitHub pull request page. The repository is tensorflow/python/keras/activations.py. A user named karmel reviewed the changes on Jul 13, 2020. The code snippet shows a diff where line 305 has been updated:

```
... @@ -302,6 +302,27 @@ def relu(x, alpha=0., max_value=None, threshold=0):  
 302     return K.relu(x, alpha=alpha, max_value=max_value, threshold=threshold)  
 303  
 304  
 305 + @keras_export('keras.activations.gelu')
```

karmel's comment: "In general, we prefer to only export to v2, as TFv1 is no longer released. Can you update this + the export below to have an empty list for v1? Eg:

```
@keras_export('keras.activations.gelu', v1=[])
```

WindQAQ's comment: "Hi @karmel, does this also apply to tf.nn.gelu, like"

```
@tf_export("nn.gelu", v1=[])
```

What is a PullRequest?

"Pull requests let you tell others about changes you've pushed to a branch in a repository on GitHub. Once a pull request is opened, you can discuss and review the potential changes with collaborators and add follow-up commits before your changes are merged into the base branch." -

<https://docs.github.com/en/pull-requests/collaborating-with-pull-requests/proposing-changes-to-your-work-with-pull-requests/about-pull-requests>



TensorFlow

THE DATA

(sample data)

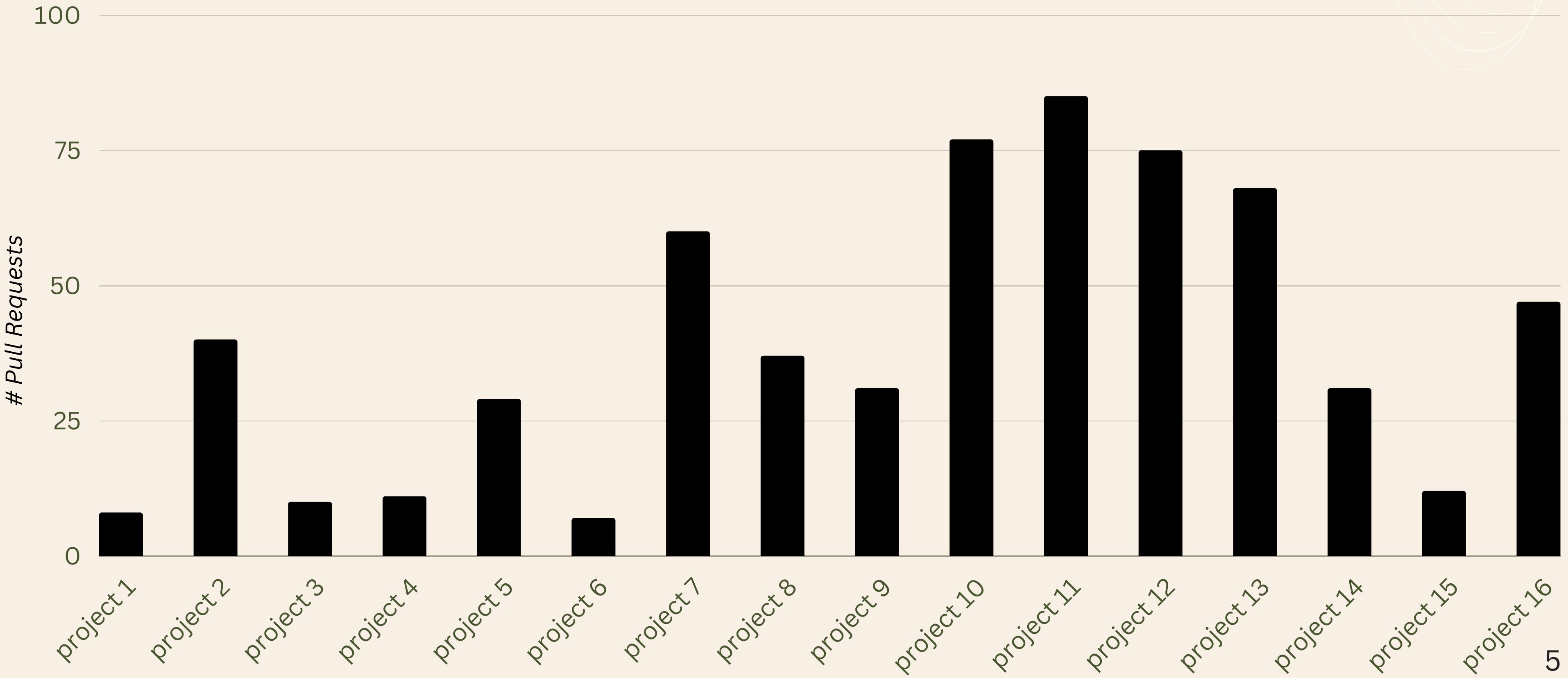


Date	Type	ID	Name	Body	URL	Comment Type
2020-07-08T01:38:31Z	Commit	8a917da9fad79e7 4772c8f3bef4ff5f8 25b88d5c	WindQAQ	COMMENTED In general we prefer to only export to v2 as TFv1 is no longer released. Can you update this + the export below to have an empty list for v1? Eg: ``` @keras_export('keras.activations.gelu' v1=[])```	https://github.com/tensorflow/tensorflow/commit/8a917da9fad79e74772c8f3bef4ff5f825b88d5c	Conventional Review

“ 16 projects in TensorFlow , each varying from 7 to 85 messages ”

THE DATA

■ Pull Requests Per Project



WHAT NOW?

2) Supervised Learning

Training a ML model with labeled data, guiding it to make accurate predictions on unseen data.

Tools: Classification Report, Heat Map

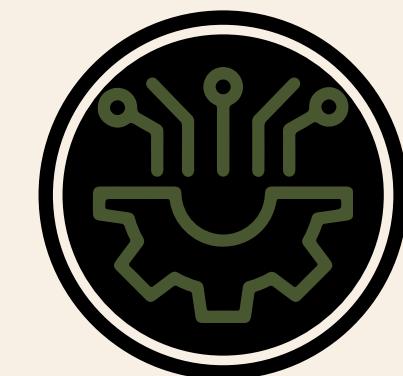
1) Manual Classification

Utilizing excel to manually represent the results of the classification.

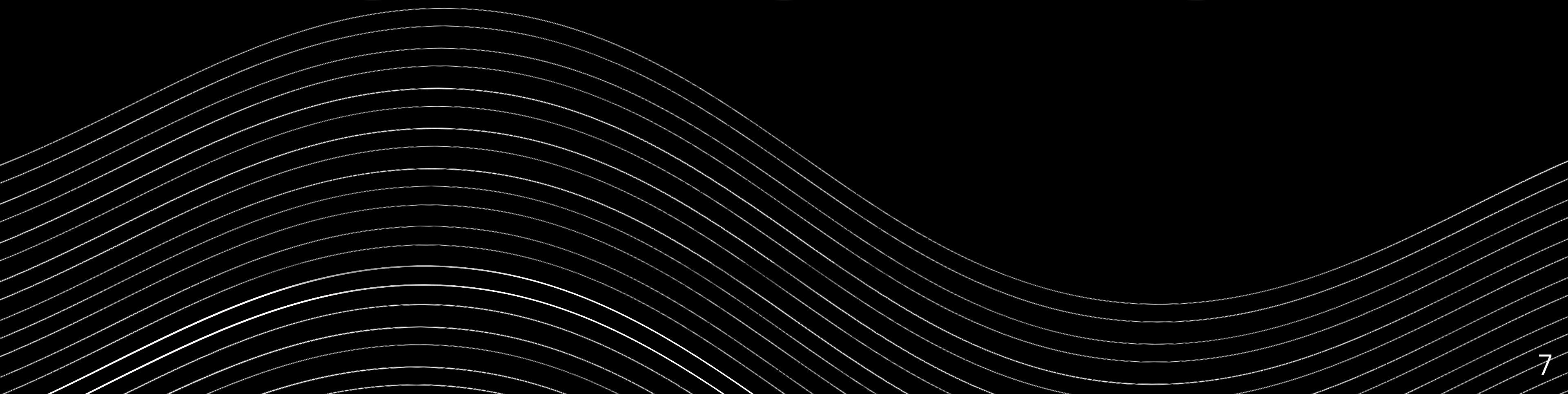
3) Unsupervised Learning

Training a ML model without labeled data, allowing it to identify patterns or structures on its own.

Tools: Gensim and the pyLDAvis library



1) MANUAL CLASSIFICATION



CATEGORIZATION PROCESS

1. Data Collection

	A	B	C	D	E	F	G	H	I	J
1	Date	Type	ID	Name	Body	URL	Comment type			
2	2020-07-08T01:38:31Z	Commit	8a917da9fad79e74772c1 WindQAQ	Migrate python implemen	https://github.com/tensorflow/tensorflow/commit/8a917da9fad79e74772c1	f83be4ff5f825b88d5c	ML review			
3	2020-07-08T01:38:42Z	Commit	9d483364d19cd7740c6e WindQAQ	Port to keras	https://github.com/tensorflow/tensorflow/commit/9d483364d19cd7740c6e93378236f1fbefb994e		ML review			
4	2020-07-08T01:41:10Z	DESCR	445836524 WindQAQ	As per https://github.com/	https://github.com/tensorflow/tensorflow/pull/41178		ML review			
5	2020-07-08T01:48:15Z	PC	655231498 seampmorgan	Thanks @WindQAQ! J	https://github.com/tensorflow/tensorflow/pull/655231498		Conventional review			
6	2020-07-08T01:55:33Z	Commit	75ad9db9f939709b5ac5d1 WindQAQ	run pylint	https://github.com/tensorflow/tensorflow/commits/75ad9db9f939709b5ac5d0a93777df8c4362b55		ML review			
7	2020-07-08T02:23:32Z	Commit	a3b21999360628de7c1 WindQAQ	Update golden api	https://github.com/tensorflow/tensorflow/commits/a3b21999360628de7c1		ML review			
8	2020-07-08T02:33:22Z	Commit	8ad98973ec85f27d9f936 WindQAQ	Update golden api	https://github.com/tensorflow/tensorflow/commits/8ad98973ec85f27d9f936		ML review			
9	2020-07-08T15:20:28Z	RC	444874482 alextx				DISMISSED			
10	2020-07-13T20:12:49Z	RC	453904111 karmel	COMMENTED Can we a	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453904111					
11	2020-07-13T20:13:06Z	RC	453904243 karmel	COMMENTED Can we a	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453904243					
12	2020-07-13T20:14:23Z	RC	447566616 karmel	CHANGES_REQUESTED!	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-447566616		Management			
13	2020-07-13T20:16:46Z	RC	453906256 karmel	COMMENTED In general	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453906256		Conventional review			
14	2020-07-13T20:33:00Z	RC	453915086 WindQAQ	COMMENTED HI @karm	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453915086		Conventional review			
15	2020-07-13T20:36:48Z	PC	657780822 WindQAQ	Hi @karmel do we also r	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-657780822		Management			
16	2020-07-14T02:39:05Z	Commit	2bcdcf38ceba1c836ed1 WindQAQ	Add example and more d	https://github.com/tensorflow/tensorflow/commits/2bcdcf38ceba1c836ed1		ML review			
17	2020-07-14T02:40:22Z	Commit	14ae0a4e8713fd92552 WindQAQ	Fix wrong module	https://github.com/tensorflow/tensorflow/pull/14ae0a4e8713fd92552		ML review			
18	2020-07-14T02:42:44Z	Commit	7d2b5cd838bacd45a8c9 WindQAQ	Export only 2 api	https://github.com/tensorflow/tensorflow/pull/7d2b5cd838bacd45a8c9		ML review			
19	2020-07-14T03:29:20Z	Commit	963717729af197d9d9c WindQAQ	Run pylint	https://github.com/tensorflow/tensorflow/pull/963717729af197d9d9d6a3d88b3bd08f5e6ca2		ML review			
20	2020-07-15T16:50:43Z	PC	658877893 hendrycks	When the C++ kernels an	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658877893		Conventional review			
21	2020-07-15T17:56:10Z	Commit	8ad3b4d1c24eb1d8805 WindQAQ	Merge branch 'master' of	https://github.com/tensorflow/tensorflow/commits/8ad3b4d1c24eb1d8805		ML review			
22	2020-07-15T20:12:27Z	RC	449294832 karmel	DISMISSED for @tensor	https://github.com/tensorflow/tensorflow/pull/449294832		Conventional review			
23	2020-07-15T20:22:43Z	PC	658989172 WindQAQ	> for @tensorflow/api-own	https://github.com/tensorflow/tensorflow/pull/658989172		Conventional review			
24	2020-07-15T20:26:57Z	PC	658991281 hendrycks	> it will somehow loss so	https://github.com/tensorflow/tensorflow/pull/658991281		Management			
25	2020-07-15T20:41:19Z	Commit	c504fe3471e214e2bf25 WindQAQ	Change approximate def	https://github.com/tensorflow/tensorflow/commits/c504fe3471e214e2bf25		ML review			
26	2020-07-15T20:44:08Z	Commit	4e69c8bafef510745c827 WindQAQ	Merge 'master' of	https://github.com/tensorflow/tensorflow/commits/4e69c8bafef510745c827		ML review			
27	2020-07-15T21:08:43Z	Commit	5a85eaab43169ab49c65 WindQAQ	Update tests	https://github.com/tensorflow/tensorflow/commits/5a85eaab43169ab49c65		ML review			
28	2020-07-15T22:22:39Z	PC	659045293 WindQAQ	Hi @karmel I change th	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-659045293		Conventional review			

2.Categorization

1

3.Visual Representation

CATEGORIZATION PROCESS

1. Data Collection

Tensor flow pull requests											
E10											
Date	Type	ID	Name	Body	URL	Comment type					
1	Commit	8a917da9fad79e74772c1 Wind/QAQ	Migrate python implemen	https://github.com/tensorflow/tensorflow/commit/8a917da9fad79e74772c1?fbclid=IwAR1f5f825b88d5c	ML review						
2	Commit	9d483364d19cd7740c6e Wind/QAQ	Port to keras	https://github.com/tensorflow/tensorflow/commit/9d483364d19cd7740c6e93378236f1fbefb994e	ML review						
3	DESCR	445836524 Wind/QAQ	As per https://github.com/	https://github.com/tensorflow/tensorflow/pull/41178	ML review						
4	PC	655231498 seampmorgan	Thanks @WindQAQ! J	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-655231498	Conventional revie						
5	Commit	75ad9db9f93709b5ac5d1 Wind/QAQ	run pylint	https://github.com/tensorflow/tensorflow/commits/75ad9db9f93709b5ac5d1#diff-093777df8c4362b55	ML review						
6	Commit	a3b21999360628de7fc1 Wind/QAQ	Update golden api	https://github.com/tensorflow/tensorflow/commits/a3b21999360628de7fc1?fbclid=IwAR533ee4a96c44fd33be	ML review						
7	Commit	8ad9c8973ec85f27d9f36 Wind/QAQ	Update golden api	https://github.com/tensorflow/tensorflow/commits/8ad9c8973ec85f27d9f36?fbclid=IwAR86164	ML review						
8	RC	444764482 aleetxp	DISMISSED	https://github.com/tensorflow/tensorflow/pull/41178#pullrequestreview-44484482	Dismissed						
9	RC	453904111 karmel	COMMENTED Can we a	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453904111	Conventional review						
10	RC	453904243 karmel	COMMENTED Can we a	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453904243	Conventional review						
11	RC	447566616 karmel	CHANGES_REQUESTED!	https://github.com/tensorflow/tensorflow/pull/41178#pullrequestreview-447566616	Management						
12	RC	453906256 karmel	COMMENTED In general	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453906256	Conventional review						
13	RC	453915086 Wind/QAQ	COMMENTED HI @karm	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453915086	Conventional review						
14	RC	657780822 Wind/QAQ	Hi @karmel do we also r	https://github.com/tensorflow/tensorflow/pull/41178#discussion_657780822	Management						
15	RC	2bcdcf38caebe1a1c836ed Wind/QAQ	Add example and more d	https://github.com/tensorflow/tensorflow/pull/41178#discussion_2bcdcf38caebe1a1c836ed	ML review						
16	Commit	14ae0a4e8713fd92552 Wind/QAQ	Fix wrong module	https://github.com/tensorflow/tensorflow/pull/41178#discussion_14ae0a4e8713fd92552	ML review						
17	Commit	7d2b5cd83b8acd45a8c9 Wind/QAQ	Export only v2 api	https://github.com/tensorflow/tensorflow/pull/41178#discussion_7d2b5cd83b8acd45a8c9	ML review						
18	Commit	963717729af197d9d9c Wind/QAQ	Run pylint	https://github.com/tensorflow/tensorflow/pull/41178#discussion_963717729af197d9d9c	ML review						
19	PC	658877893 hendrycks	When the C++ kernels an	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658877893	Conventional revie						
20	Commit	8ad3b4d1c24eb1d8805 Wind/QAQ	Merge branch 'master' of	https://github.com/tensorflow/tensorflow/pull/41178#pullrequestreview-658877894	ML review						
21	Commit	449294833 karmel	DISMISSED for @tensord	https://github.com/tensorflow/tensorflow/pull/41178#pullrequestreview-449294833	Conventional revie						
22	RC	658989121 Wind/QAQ	> for @tensorflow/api-own	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658989121	Management						
23	PC	658991281 hendrycks	> it will somehow loss so	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658991281	ML review						
24	PC	c504fe3471e214e2bf25 Wind/QAQ	Change approximate def	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658991281	ML review						
25	Commit	4669c8bafef510745c827 Wind/QAQ	Merge 'master' of	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658994287	ML review						
26	Commit	5a85ea8b4316984b9c65 Wind/QAQ	Update tests	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-659045293	ML review						
27	PC	659045283 Wind/QAQ	Hi @karmel I change th	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-659045293	Conventional revie						
28	Commit										

Comment type	Nature of comments	Contributed artifacts	Why	Relevant SDLC or ML phase
ML review	Clarification	Code	Code changes	SWEBOK-Software engineering process
ML review	Clarification	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	other	Assisting the changes	SWEBOK-Software engineering process
Conventional review	Comment	Review	Question	SWEBOK-Software engineering management
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
Dismissed	Dismissed	Dismissed	Dismissed	Dismissed
Conventional review	Comment	Link	Asked question	SWEBOK-Software engineering management
Conventional review	Comment	Code	Asked question	SWEBOK-Software engineering management
Management	Pinging	other	Process management	SWEBOK-Software engineering management
Conventional review	Clarification	Code	Asked question	SWEBOK-Software engineering process
Conventional review	Comment	Code	Question	SWEBOK-Software engineering management
Management	Clarification	Code	Assisting the changes	SWEBOK-Software engineering management
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Feedback	Code	Explaining changes	SWEBOK-Software engineering management
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
Conventional review	Comment	other	Explaining plan	SWEBOK-Software engineering management
Conventional review	Comment	other	Code review	SWEBOK-Software engineering management
Management	Code review	other	Assisting the changes	SWEBOK-Software engineering management
ML review	Code review	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code review	Code	Approved changes	SWEBOK-Software engineering process
ML review	Code review	Code	Approved changes	SWEBOK-Software engineering management
Conventional review	Request a review	Review	Code review	SWEBOK-Software engineering management

3. Visual Representation

CATEGORIZATION PROCESS

1. Data Collection

Tensor flow pull requests											
E10											
Date	Type	ID	Name	Body	URL	Comment type	Nature of comments	Contributed artifacts	Why	Relevant SDLC or ML phase	
2020-07-08T01:38:31Z	Commit	8a917da9fad79e74772c1 WindQAQ	Migrate python implemen	https://github.com/tensorflow/tensorflow/commit/8a917da9fad79e74772c1?fbclid=IwAR1f825b88d5c	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Clarification	Code	Code changes	SWEBOK-Software engineering process	
2020-07-08T01:38:42Z	Commit	9d483364d19cd7740c6e WindQAQ	Port to keras	https://github.com/tensorflow/tensorflow/commit/9d483364d19cd7740c6e93378236f1fbefb994e	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Clarification	Code	Assisting the changes	SWEBOK-Software engineering process	
2020-07-08T01:41:10Z	DESCR	445836524 WindQAQ	As per https://github.com/	https://github.com/tensorflow/tensorflow/pull/41178	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Code	other	Assisting the changes	SWEBOK-Software engineering process	
2020-07-08T01:48:15Z	PC	655231498 seampongman	Thanks @WindQAQ! J	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-655231498	https://github.com/tensorflow/tensorflow/pull/41178	Conventional revie	Comment	Review	Question	SWEBOK-Software engineering management	
2020-07-08T01:55:33Z	Commit	75ad9db9f93709b5ac5d1 WindQAQ	run pylint	https://github.com/tensorflow/tensorflow/commit/75ad9db9f93709b5ac5d0a93777df8c4362b55	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process	
2020-07-08T02:33:32Z	Commit	a3b21999360628de7fc1 WindQAQ	Update golden api	https://github.com/tensorflow/tensorflow/commit/a3b21999360628de7fc1?fbclid=IwAR533ee4a96c44fd33be	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process	
2020-07-08T15:20:28Z	Commit	8ad9c8973ec85f27d9f36 WindQAQ	Update golden api	https://github.com/tensorflow/tensorflow/commit/8ad9c8973ec85f27d9f36	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process	
2020-07-08T15:20:28Z	RC	444764482 alextp	DISMISSED	https://github.com/tensorflow/tensorflow/pull/41178#pullrequestreview-44484482	https://github.com/tensorflow/tensorflow/pull/41178	Dismissed	Comment	Link	Asked question	SWEBOK-Software engineering management	
2020-07-13T20:12:49Z	RC	453904111 karmel	COMMENTED Can we a	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453904111	https://github.com/tensorflow/tensorflow/pull/41178	Conventional revie	Comment	Code	Asked question	SWEBOK-Software engineering management	
2020-07-13T20:13:06Z	RC	453904243 karmel	COMMENTED Can we a	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453904243	https://github.com/tensorflow/tensorflow/pull/41178	Conventional revie	Comment	Code	Asked question	SWEBOK-Software engineering management	
2020-07-13T20:14:23Z	RC	447566616 karmel	CHANGES_REQUESTER	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-447566616	https://github.com/tensorflow/tensorflow/pull/41178	Management	Pinging	other	Process management	SWEBOK-Software engineering management	
2020-07-13T20:16:46Z	RC	453906256 karmel	COMMENTED In general	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453906256	https://github.com/tensorflow/tensorflow/pull/41178	Conventional revie	Clarification	Code	Asked question	SWEBOK-Software engineering process	
2020-07-13T20:33:00Z	RC	453915086 WindQAQ	COMMENTED HI @karmel	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453915086	https://github.com/tensorflow/tensorflow/pull/41178	Conventional revie	Comment	Code	Question	SWEBOK-Software engineering management	
2020-07-13T20:36:48Z	RC	657780822 WindQAQ	Hi @karmel do we also r	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-657780822	https://github.com/tensorflow/tensorflow/pull/41178	Management	Clarification	Code	Assisting the changes	SWEBOK-Software engineering management	
2020-07-14T02:39:05Z	Commit	2bcdcf38caebe1a836ed WindQAQ	Add example and more d	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453906256	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Comment	Code	Assisting the changes	SWEBOK-Software engineering process	
2020-07-14T02:40:22Z	Commit	14ae0a4e8713fd92552 WindQAQ	Fix wrong module	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453906256	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Comment	Code	Assisting the changes	SWEBOK-Software engineering process	
2020-07-14T02:42:42Z	Commit	7d2b5cd838bacd45a8c9 WindQAQ	Export only v2 api	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453906256	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Comment	Code	Assisting the changes	SWEBOK-Software engineering process	
2020-07-14T03:29:20Z	Commit	963717729af197d9d9c WindQAQ	Run pylint	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453906256	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process	
2020-07-15T16:50:43Z	PC	658877893 hendrycks	When the C++ kernels an	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658877893	https://github.com/tensorflow/tensorflow/pull/41178	Conventional revie	Feedback	Code	Explaining changes	SWEBOK-Software engineering management	
2020-07-15T17:56:10Z	Commit	8ad3b4d1c24eb1d8805 WindQAQ	Merge branch 'master' of	https://github.com/tensorflow/tensorflow/pull/41178#discussion_453906256	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Comment	Code	Explaining changes	SWEBOK-Software engineering management	
2020-07-15T20:12:27Z	RC	44924833 karmel	DISMISSED for @tensord	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-44924833	https://github.com/tensorflow/tensorflow/pull/41178	Conventional revie	Comment	Code	Assisting the changes	SWEBOK-Software engineering process	
2020-07-15T22:22:43Z	PC	658989172 WindQAQ	> for > will somehow loss so	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658989172	https://github.com/tensorflow/tensorflow/pull/41178	Management	Comment	other	Explaining plan	SWEBOK-Software engineering management	
2020-07-15T20:26:57Z	PC	658991281 hendrycks	Change approximate def	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658991281	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Comment	other	Code review	SWEBOK-Software engineering management	
2020-07-15T20:41:19Z	Commit	c504fe3471e24bf25 WindQAQ	ML review	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658991281	https://github.com/tensorflow/tensorflow/pull/41178	Management	Comment	other	Code review	SWEBOK-Software engineering management	
2020-07-15T20:44:08Z	Commit	4e69c8baf5f10745c827 WindQAQ	Merge branch 'master' of	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-658991281	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Comment	other	Code review	SWEBOK-Software engineering management	
2020-07-15T21:08:43Z	Commit	5a85ea8b4316984b9c65 WindQAQ	Update tests	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-659045293	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Code review	Code	Approved changes	SWEBOK-Software engineering process	
2020-07-15T22:22:39Z	PC	659045283 WindQAQ	Hi @karmel I change th	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-659045293	https://github.com/tensorflow/tensorflow/pull/41178	ML review	Code review	Code	Approved changes	SWEBOK-Software engineering management	

2.Categorization

Comment type	Nature of comments	Contributed artifacts	Why	Relevant SDLC or ML phase
ML review	Clarification	Code	Code changes	SWEBOK-Software engineering process
ML review	Clarification	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	other	Assisting the changes	SWEBOK-Software engineering process
Conventional review	Comment	Review	Question	SWEBOK-Software engineering management
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
Dismissed	Dismissed	Dismissed	Dismissed	Dismissed
Conventional review	Comment	Link	Asked question	SWEBOK-Software engineering management
Conventional review	Comment	Code	Asked question	SWEBOK-Software engineering management
Management	Pinging	other	Process management	SWEBOK-Software engineering management
Conventional review	Clarification			

CATEGORIZATION PROCESS

1. Data Collection

Tensor flow pull requests											
File Edit View Insert Format Data Tools Extensions Help											
E10											
1 Date											
2	2020-07-08T01:38:31Z	Type	ID	Name	Body	URL	Comment type				
3	Commit	8a917da9fad79e74772c1	WindQAQ	Migrate python implemen	https://github.com/tensorflow/tensorflow/commit/8a917da9fad79e74772c1ff3be4ff5f825b88d5c	ML review					
4	Commit	9d483364d19cd7740c6e	WindQAQ	Port to keras	https://github.com/tensorflow/tensorflow/commit/9d483364d19cd7740c6e93378236f1f9febf994e	ML review					
5	2020-07-08T01:41:10Z	DESCR	445836524	WindQAQ	As per https://github.com/	https://github.com/tensorflow/tensorflow/pull/41178	ML review				
6	PC	655231498	seanmpmorgan	Thanks @WindQAQ! J	https://github.com/tensorflow/tensorflow/pull/41178#issuecomment-655231498	Conventional revie					
7	2020-07-08T01:55:33Z	Commit	75ad9db9f939709b5ac5d1	WindQAQ	run pylint	https://github.com/tensorflow/tensorflow/commits/75ad9db9f939709b5ac5d1f8c4362b55	ML review				
8	2020-07-08T03:23:32Z	Commit	a3b21999360628de7fc1	WindQAQ	Update golden api	https://github.com/tensorflow/tensorflow/commits/a3b21999360628de7fc1e553eeea96e44fd33be	ML review				
9	2020-07-08T15:20:28Z	Commit	8ad9c8973e85f27d9f36	WindQAQ	Update golden api	https://github.com/tensorflow/tensorflow/commits/8ad9c8973e85f27d9f362abeb3020270a86164	ML review				
10	2020-07-08T15:20:49Z	RC	444764482	alextp	DISMISSED	https://github.com/tensorflow/tensorflow/pull/41178/pullrequestreview-444874482	Dismissed				
11	2020-07-13T20:13:06Z	RC	453904243	karmel	COMMENTED Can we a	https://github.com/tensorflow/tensorflow/pull/41178/disscussion_453904243	Conventional revie				
12	2020-07-13T20:14:23Z	RC	447566616	karmel	CHANGES_REQUESTE	https://github.com/tensorflow/tensorflow/pull/41178/pullrequestreview-447566616	Management				
13	2020-07-13T20:16:46Z	RC	453906256	karmel	COMMENTED In general	https://github.com/tensorflow/tensorflow/pull/41178/disscussion_453906256	Conventional revie				
14	2020-07-13T20:33:00Z	RC	453915086	WindQAQ	COMMENTED Hi @karm	https://github.com/tensorflow/tensorflow/pull/41178/disscussion_453915086	Conventional revie				
15	2020-07-13T20:36:48Z	PC	657780822	WindQAQ	Hi @karmel do we also r	https://github.com/tensorflow/tensorflow/pull/41178/disscussion_657780822	Management				
16	2020-07-14T02:39:05Z	Commit	2bcdcf38caebe1a836ed	WindQAQ	Add example and more d	https://github.com/tensorflow/tensorflow/commit/2bcdcf38caebe1a836ed02d0da5e451ce824275	ML review				
17	2020-07-14T02:40:22Z	Commit	14aae0a4e8713fd92552	WindQAQ	Fix wrong module	https://github.com/tensorflow/tensorflow/commit/14aae0a4e8713fd92552d98ff4127092c4b7b81	ML review				
18	2020-07-14T02:42:44Z	Commit	7d2b5cd83b8acd45a8c9	WindQAQ	Export only v2 api	https://github.com/tensorflow/tensorflow/commit/7d2b5cd83b8acd45a8c93309a5d01164ace5cd	ML review				
19	2020-07-14T03:29:20Z	Commit	963717729af197d9df9c	WindQAQ	Run pylint	https://github.com/tensorflow/tensorflow/commit/963717729af197d9df9d6a3d88b3bd085e6ca2	ML review				
20	2020-07-15T16:50:43Z	PC	658877893	hendrycks	When the C++ kernels an	https://github.com/tensorflow/tensorflow/pull/41178/issucomment-658877893	Conventional revie				
21	2020-07-15T16:56:10Z	Commit	8ad3b4d1c24ebe1d8805	WindQAQ	Merge branch 'master' of	https://github.com/tensorflow/tensorflow/commit/8ad3b4d1c24ebe1d8805ced6b7663a4228694	ML review				
22	2020-07-15T20:12:27Z	RC	449294833	karmel	DISMISSED for @tens	https://github.com/tensorflow/tensorflow/pull/41178/pullrequestreview-449294833	Conventional revie				
23	2020-07-15T20:22:43Z	PC	658989172	WindQAQ	> for @tensorflow/api-own	https://github.com/tensorflow/tensorflow/pull/41178/issucomment-658989172	Conventional revie				
24	2020-07-15T20:26:57Z	PC	658991281	hendrycks	it will somehow loss so	https://github.com/tensorflow/tensorflow/pull/41178/issucomment-658991281	Management				
25	2020-07-15T20:41:19Z	Commit	c504fee3471e214e2bf25	WindQAQ	Change approximate def	https://github.com/tensorflow/tensorflow/commit/c504fee3471e214e2bf25fb757894db90b4395a5b5	ML review				
26	2020-07-15T20:44:08Z	Commit	4e69cb8afe5f10745c827	WindQAQ	Merge branch 'master' of	https://github.com/tensorflow/tensorflow/commit/4e69cb8afe5f10745c8279c7d712d2df23beb	ML review				
27	2020-07-15T21:08:43Z	Commit	5a85ea8b4316984b9c65	WindQAQ	Update tests	https://github.com/tensorflow/tensorflow/commit/5a85ea8b4316984b9c6581d81fa002a345fb8af	ML review				
28	2020-07-15T22:22:39Z	PC	659045283	WindQAQ	Hi @karmel I change th	https://github.com/tensorflow/tensorflow/pull/41178/issucomment-659045283	Conventional revie				

Comment type	Nature of comments	Contributed artifacts	Why	Relevant SDLC or ML phase
ML review	Clarification	Code	Code changes	SWEBOK-Software engineering process
ML review	Clarification	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	other	Assisting the changes	SWEBOK-Software engineering process
Conventional review	Comment	Review	Question	SWEBOK-Software engineering management
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
Dismissed	Dismissed	Dismissed	Dismissed	Dismissed
Conventional review	Comment	Link	Asked question	SWEBOK-Software engineering management
Conventional review	Comment	Code	Asked question	SWEBOK-Software engineering management
Management	Pinging	other	Process management	SWEBOK-Software engineering management
Conventional review	Clarification	Code	Asked question	SWEBOK-Software engineering process
Conventional review	Comment	Code	Question	SWEBOK-Software engineering management
Management	Clarification	Code	Assisting the changes	SWEBOK-Software engineering management
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Feedback	Code	Explaining changes	SWEBOK-Software engineering management
ML review	Code	Code	Explaining changes	SWEBOK-Software engineering process
ML review	Code	Code	Assisting the changes	SWEBOK-Software engineering process
Conventional review	Comment	other	Explaining plan	SWEBOK-Software engineering management
Conventional review	Comment	other	Code review	SWEBOK-Software engineering management
Management	Code review	other	Assisting the changes	SWEBOK-Software engineering management
ML review	Code review	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code review	Code	Assisting the changes	SWEBOK-Software engineering process
ML review	Code review	Code	Approved changes	SWEBOK-Software engineering process
ML review	Code review	Code	Approved changes	SWEBOK-Software engineering management
Conventional review	Request a review	Review	Code review	SWEBOK-Software engineering management

6	ML review
7	ML review
8	Conventional review
9	Management
10	Other
11	Management
12	Management
13	Management
14	Management
15	Management

<u

Manual Classification

PROS

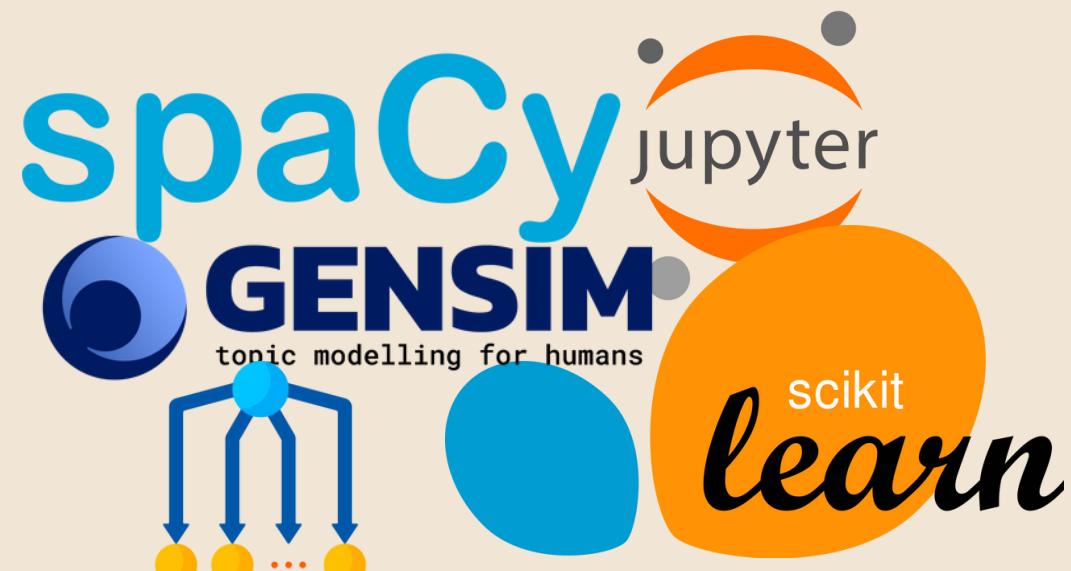
- Allowed Data Familiarization
- Quality Control
 - Able to double check and compare with other results
- Precision and Accuracy
 - Human judgement combined with collaboration

CONS

- Time-Consuming & Costly
- Oversight/ Inconsistencies
- Human Error: Typos, Bias

HOW CAN WE IMPROVE?

INTRODUCING Machine learning TO THE PROJECT...



GETTING READY TO TRAIN!

- Code Breakdown:
 - a. Import a .csv file containing a combined version of all the sheets
 - b. Give a Category Value for each sub-category in "Comment type" column
 - c. Remove stop words and lemmatize the text from the "Body" column, store them in a vector list
 - d. Give each vector a vectorized value (turning text into a numerical representation)
 - e. Create training and test values

```
#Read in and print the csv file
df = pd.read_csv("combined_tf_pr.csv")
print(df.shape)
#print the 39 rows
df.head(607)
```

GETTING READY TO TRAIN!

- Code Breakdown:
 - a.~~Import a .csv file containing a combined version of all the sheets~~
 - b. **Give a Category Value for each sub-category in "Comment type" column**
 - c. Remove stop words and lemmatize the text from the "Body" column, store them in a vector list
 - d. Give each vector a vectorized value (turning text into a numerical representation)
 - e. Create training and test values

Why is this important?

Ensures efficiency in data processing!

```
#New Column: Gives a unique number to each of these categories
df['category_val'] = df['Comment type'].map({'Conventional review' : 0, 'ML review': 1, 'Dismissed' : 2, 'Management' : 3, 'Other' : 4, 'initial implementation' : 5, 'empty' : 6})
df.head(607)
```

GETTING READY TO TRAIN!

- Code Breakdown:
 - a. Import a .csv file containing a combined version of all the sheets
 - b. Give a Category Value for each sub-category in "Comment type" column
 - c. Remove stop words and lemmatize the text from the "Body" column, store them in a vector list
 - d. Give each vector a vectorized value (turning text into a numerical representation)
 - e. Create training and test values

***Why use nltk
and spacy?***

Ensures better coverage!

```
print(nltk_stopwords)
print("\n")
print(spacy_stopwords)
print("\n")
print(combined_stopwords)

['i', 'me', 'my', 'myself', 'we', 'our', 'ours', 'ourselves', 'you', "you're", "you've", "you'll", "you'd", 'your', 'yours', 'yourself', 'yo

{'side', 'due', 'might', 'unless', 'hereby', 'us', 'a', 'towards', 'forty', 'thereafter', 'have', 'beforehand', 'this', 'his', 'no', 'seemed

{'side', 'hereupon', 'y', 'out', 'aren', 'does', 'due', 'before', "shan't", 'might', 'done', 'namely', 'unless', 'hereby', 'us', 'how', 'ele
```

GETTING READY TO TRAIN!

- Code Breakdown:
 - a. Import a .csv file containing a combined version of all the sheets
 - b. Give a Category Value for each sub-category in "Comment type" column
 - c. Remove stop words and lemmatize the text from the "Body" column, store them in a vector list
 - d. Give each vector a vectorized value (turning text into a numerical representation)
 - e. Create training and test values

```
#New Column: Gives. vectorized value of each input value
df['vector'] = df['Body'].apply(lambda text: preprocess_and_vectorize(text))
df.head(607)
```

```
[array([-2.4752 ,  0.52168 ,  0.16867 ,  3.2994 ,  2.1103 ,  0.25 ,
       0.046571,  2.7635 , -1.3151 ,  1.1516 ,  3.4943 ,  1.0532 ,
      -4.1382 ,  3.3488 ,  2.7279 ,  0.023003,  1.3536 ,  1.4741 ,
      -2.6494 , -2.9689 , -2.6738 , -2.7591 , -1.2577 ,  1.3407 ,
       0.58053 ,  0.32727 , -3.3395 , -0.35698 , -1.2736 ,  1.3963 ,
      -1.476 , -1.4323 , -1.1268 , -0.25453 , -1.9975 , -2.7572 ,
       1.3462 , -1.0832 ,  0.91696 , -3.0349 , -0.76496 , -1.9992 ,
       0.063289,  3.5607 , -2.7648 ,  0.76119 ,  2.2463 , -1.5844 ,
      -2.0389 , -1.5829 ,  1.0624 ,  2.8675 , -2.8235 , -4.5756 ,
      -2.7501 , -0.085979, -5.9145 ,  1.1754 ,  1.1318 , -0.38072 ,
      -3.2722 , -2.1684 , -0.4984 ,  0.88358 ,  5.0679 ,  2.0915 ,
      -2.7009 ,  0.22031 ,  2.2098 , -0.36808 ,  2.1906 , -3.6019 ,
       4.1686 ,  1.5225 , -0.43594 ,  0.42891 , -1.6773 , -0.28204 ,
      -0.59196 ,  1.0297 , -2.7454 , -0.89824 , -2.0003 ,  0.72807 ,
       1.0572 , -0.78476 ,  1.3047 , -4.2398 , -0.47198 ,  0.21958 ,
       2.6356 ,  0.53105 , -2.7675 ,  0.1523 , -0.71488 ,  2.7661 ,
       0.22556 , -0.48985 , -0.24394 ,  0.26594 ,  2.0607 ,  0.61695 ,
```

Why is this important?

helps bridge the gap between
human and machine.

eg.

woman -> queen, man -> king

GETTING READY TO TRAIN!

- Code Breakdown:
 - a. Import a .csv file containing a combined version of all the sheets
 - b. Give a Category Value for each sub-category in "Comment type" column
 - c. Remove stop words and lemmatize the text from the "Body" column, store them in a vector list
 - d. Give each vector a vectorized value (turning text into a numerical representation)
 - e. Create training and test values

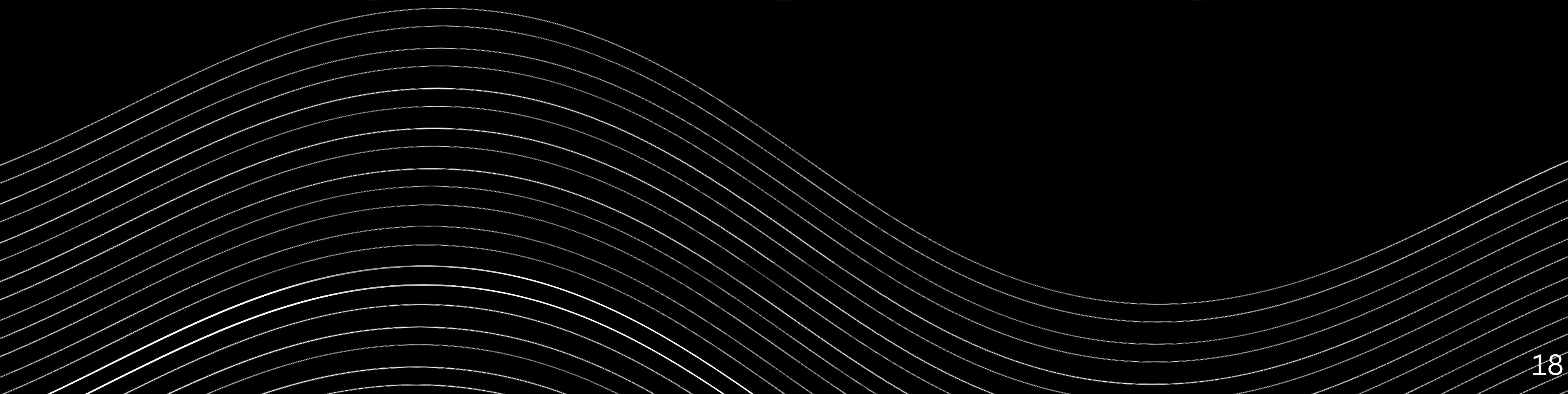
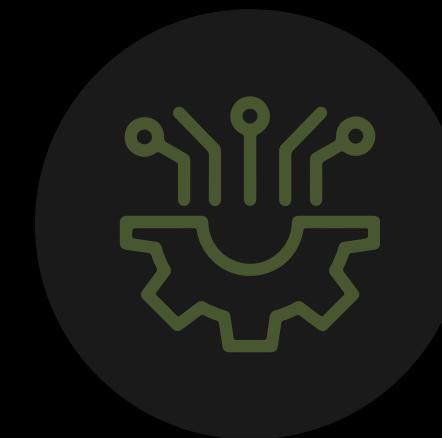
Why is this important?

Allows the model to learn patterns from the training data and then its accuracy can be assessed on unseen data,

```
#Splitting the vectors of df into sets of X_train and testing with X_test. then splitting category_val into Y_train and testing with Y_test
X_train, X_test, Y_train, Y_test = train_test_split(
    df['vector'],    # The data needing to split
    df['category_val'], # the category_values associated to each vector
    test_size=0.2,   #testing 20% and training 80% of the data set
    random_state =2023 #ensuring consistency in splitting
)
```

training (80%),
testing (20%)

SUPERVISED LEARNING



CLASSIFICATION REPORT

Precision: The proportion of positive identifications that were actually correct. A model that produces no false positives has a precision of 1.0.

Recall: The proportion of actual positives that were identified correctly.

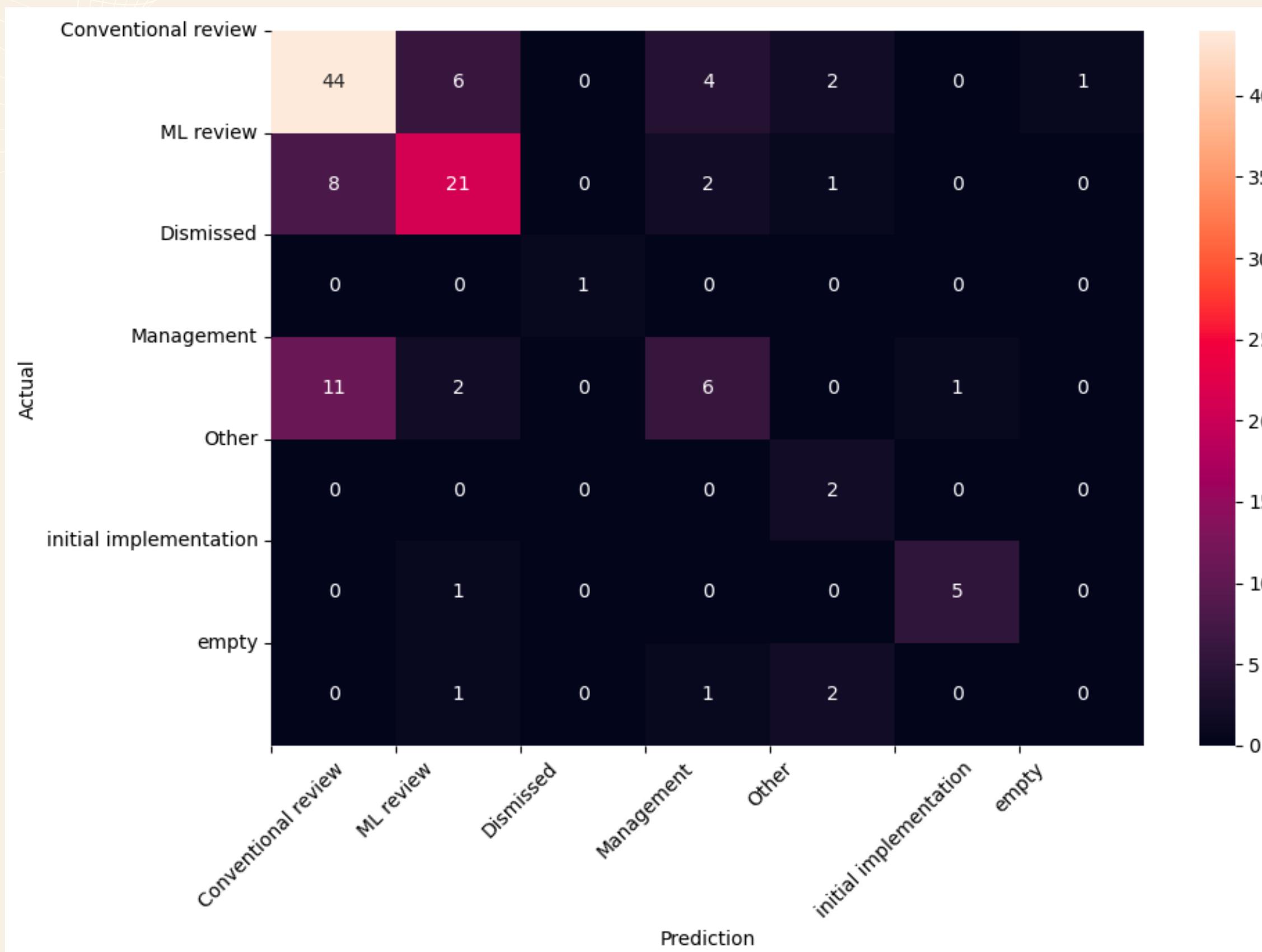
f1-score: Avg of precision and recall value. Best value at 1 and worst at 0.

Accuracy: Proportion of all predictions that were correct.

	labels	precision	recall	f1-score
Conventional review	0	0.70	0.77	0.73
ML review	1	0.68	0.66	0.67
Dismissed	2	1.00	1.00	1.00
Management	3	0.46	0.30	0.36
Other	4	0.29	1.00	0.44
initial implementation	5	0.83	0.83	0.83
empty	6	0.00	0.00	0.00
accuracy				0.65

*The overall accuracy is 0.65
-> the model correctly predicted the class 65% of the time on the test set!*

CONFUSION MATRIX



What does this tell us? Identifies the overall accuracy of the model AND also which specific classes the model is struggling with.

How to read?

- The diagonal from the top-left to the bottom-right represent the true positive counts for each class
- The rest are false predictions
- A higher TP values = higher accuracy.
- Actual:
- Prediction:

Supervised Learning

PROS

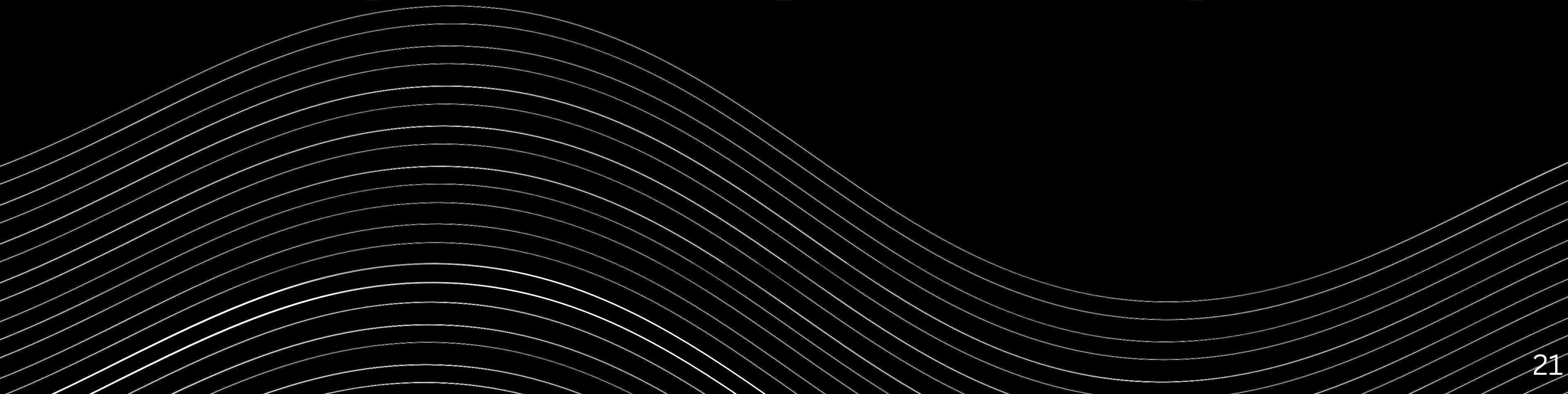
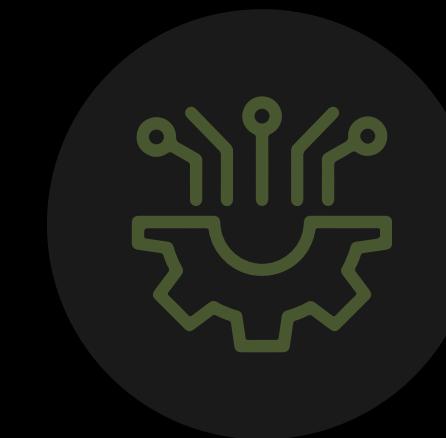
- High Accuracy
 - Assumption: The Given Data is Also Accurate
- Strong Predictive Power
- Data and Time Efficiency

CONS

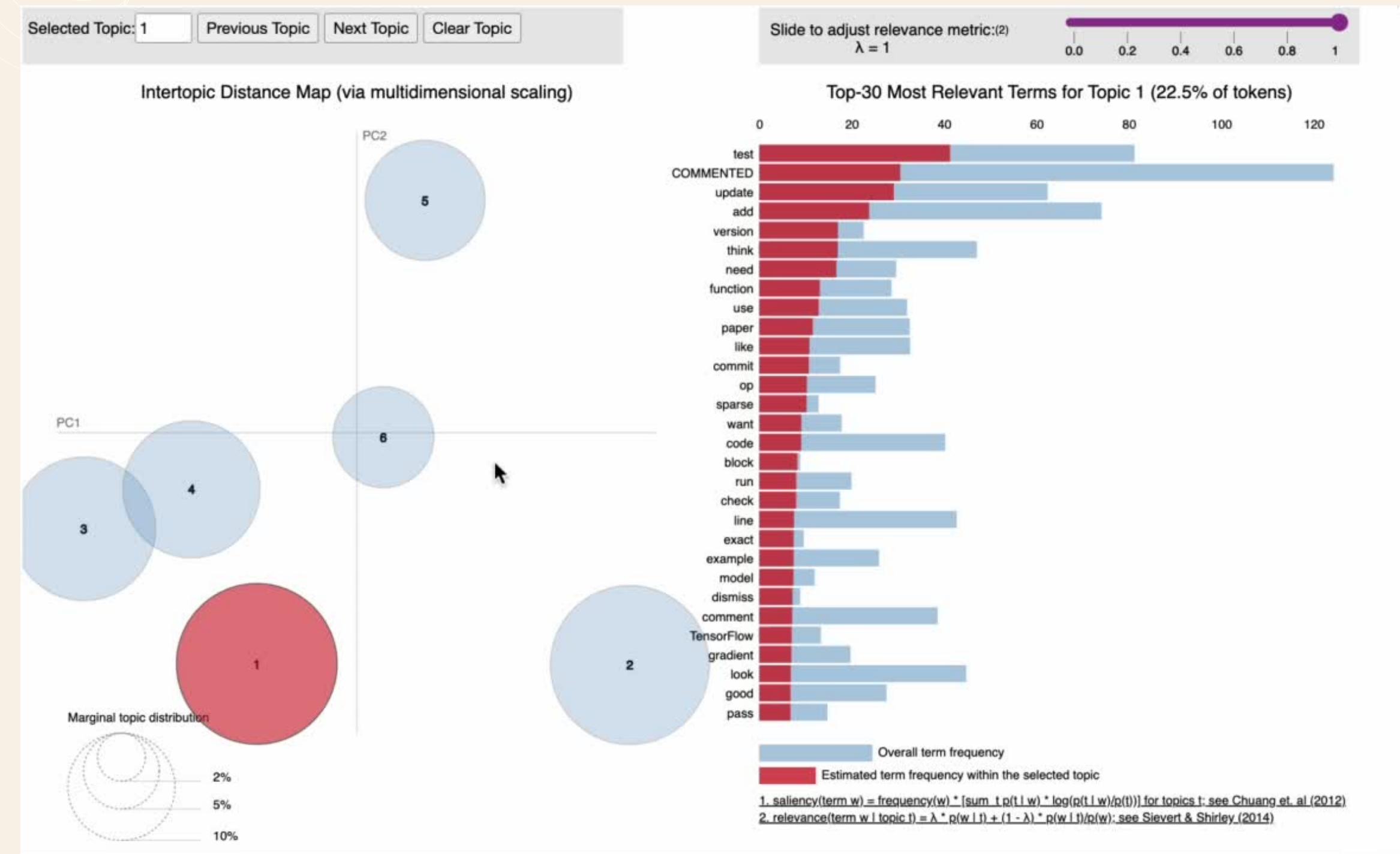
- Need to Provide Labeled Data: Time-Consuming & Costly Process
- Oversight/ Inconsistencies: From Bias results That The Machine Can Pick Up On
- Data Processing:
 - Lots of Time is Needed to Review the Results and Ensure Accuracy

HOW CAN WE IMPROVE?

UNSUPERVISED LEARNING



PYLDAVIS GENSIM



What does this tell us?

- How many topics there are and their relevance.
- How topics relate to each other.
- The most important terms in the entire corpus and within each topic.

How to read?

- On the Left:
 - The larger the bubble the bigger the topic relevance
 - Distance between bubbles indicates the similarity between topics – closer together are more similar
- On the right:
 - This list shows the overall most descriptive/relevant terms in the corpus.

Unsupervised Learning

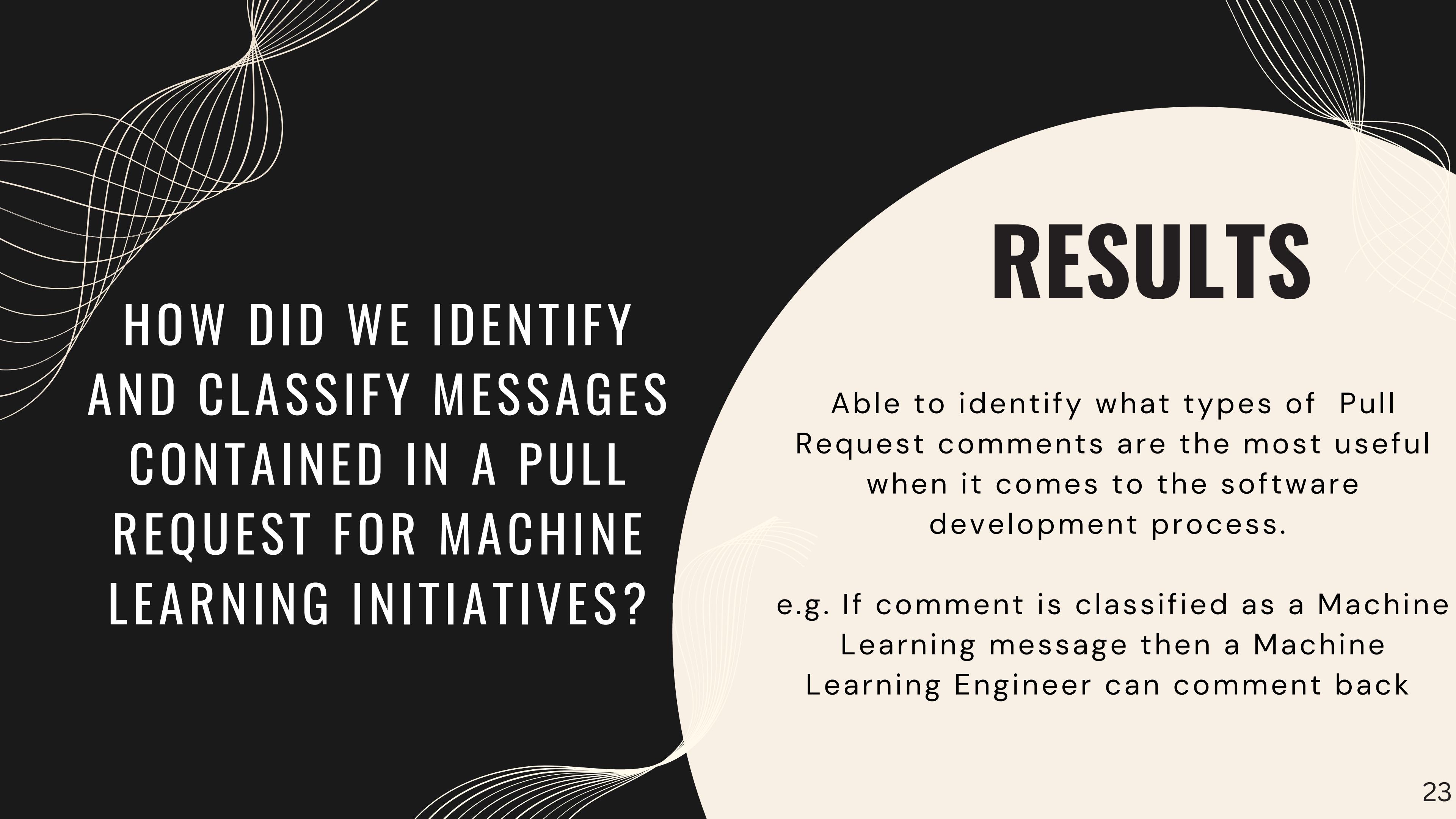
PROS

- No Need for Labeled Data
- Great for Data Exploration
- Data and Time Efficiency

CONS

- No clear objective with the data
- Sensitivity for Preprocessing
- Data Processing:
 - Lots of Time is Needed to Review the Results and Ensure Accuracy

HOW CAN WE IMPROVE?



HOW DID WE IDENTIFY AND CLASSIFY MESSAGES CONTAINED IN A PULL REQUEST FOR MACHINE LEARNING INITIATIVES?

RESULTS

Able to identify what types of Pull Request comments are the most useful when it comes to the software development process.

e.g. If comment is classified as a Machine Learning message then a Machine Learning Engineer can comment back

FUTURE STEPS

Expand the code coverage to include the rest of the Categories such as "Nature of Comments", "Contributed Artifacts" etc..

BUMPS/OBSTACLES

Understanding and scheduling the amount of time necessary to complete each part of the project.

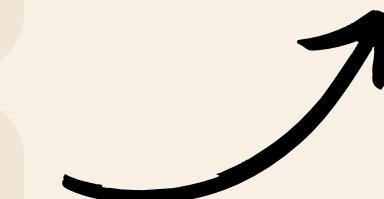
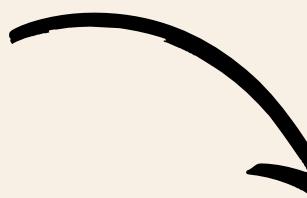
TIME
MANAGEMENT

Reading and researching the best ways to approach the machine learning aspect of the project based on my current limited knowledge

ML KNOWLEDGE

Understanding the various outputs and reflecting on what they mean based on the given data and how to fix the issues that are being represented incorrectly

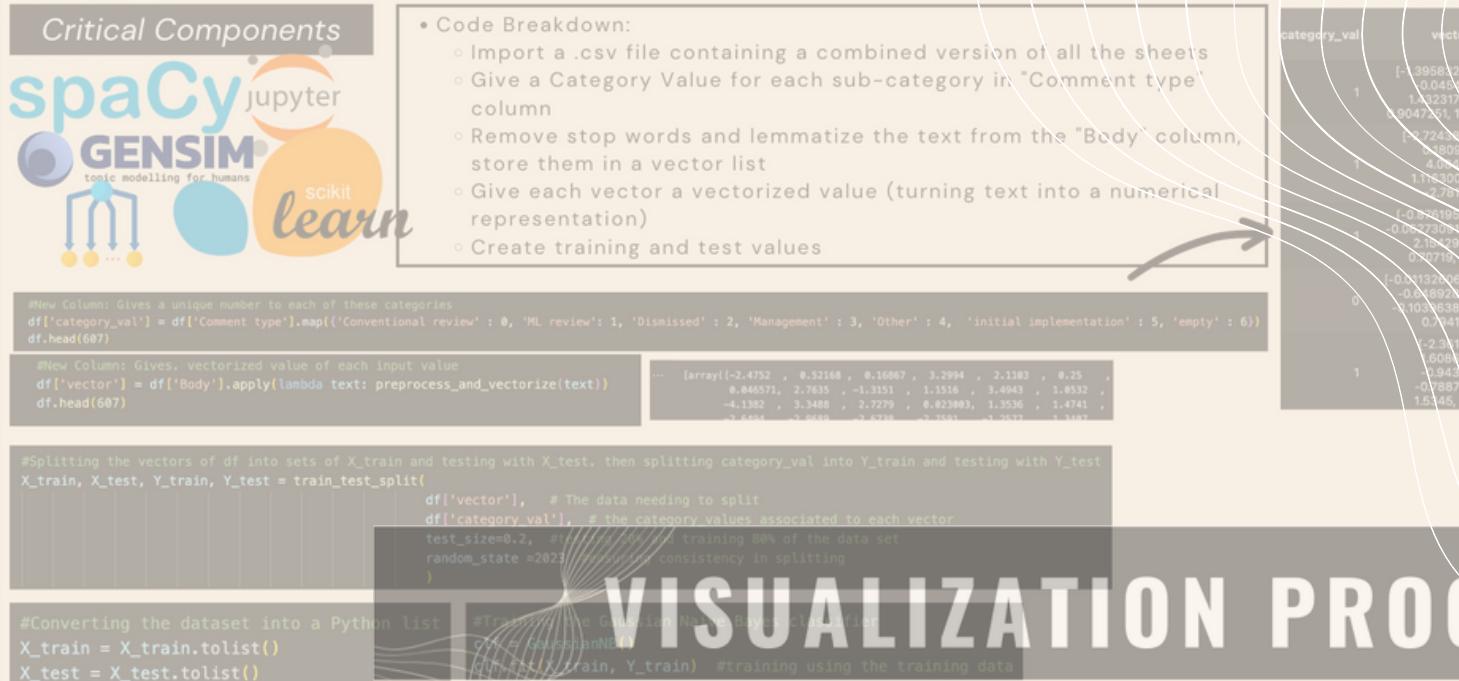
DEBUGGING



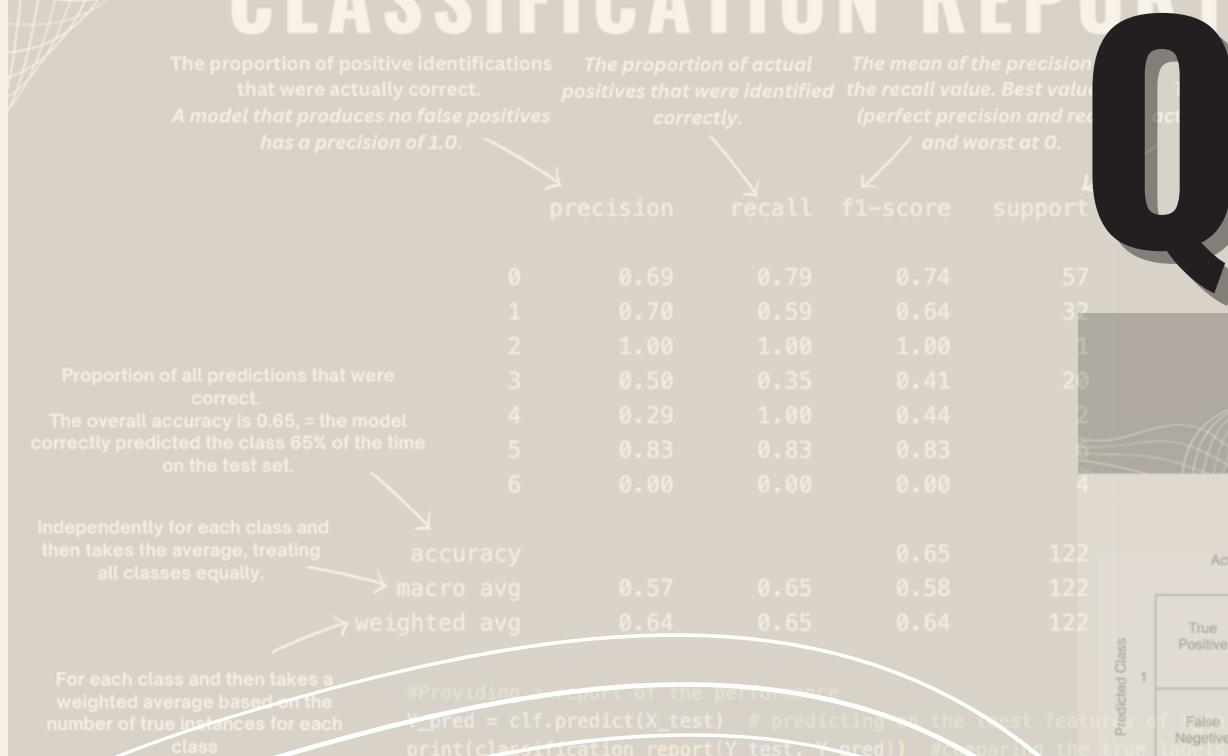
CATEGORIZATION PROCESS



IMPLEMENTATION PROCESS

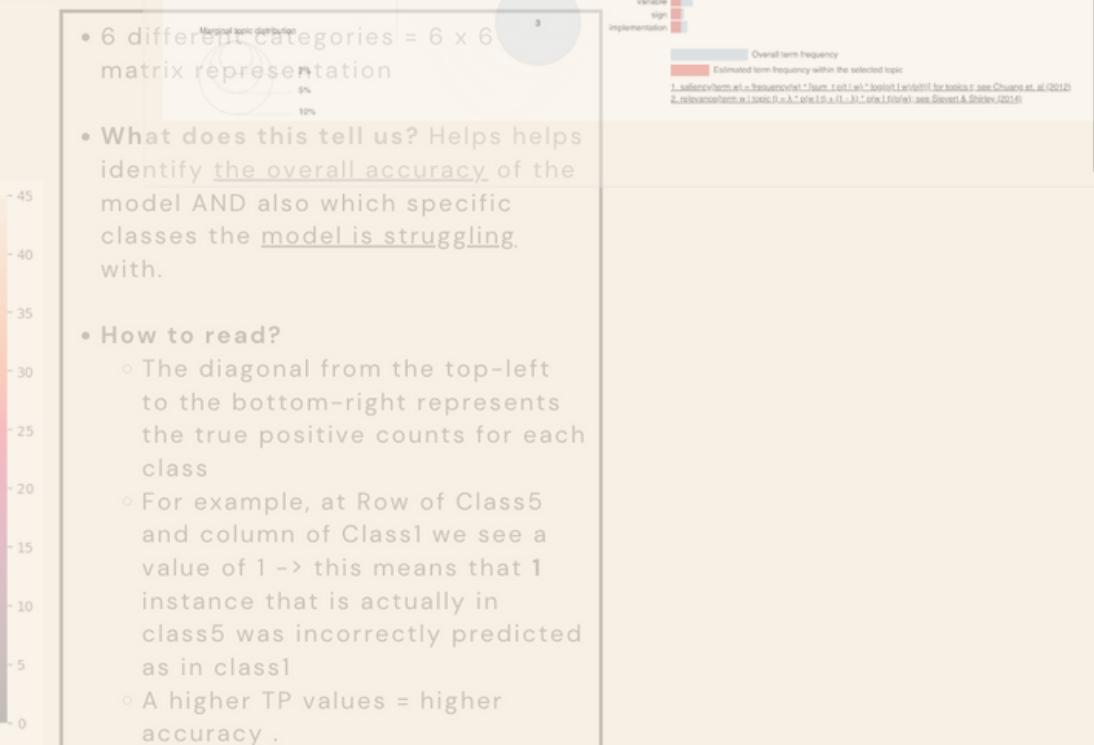


CLASSIFICATION REPORT

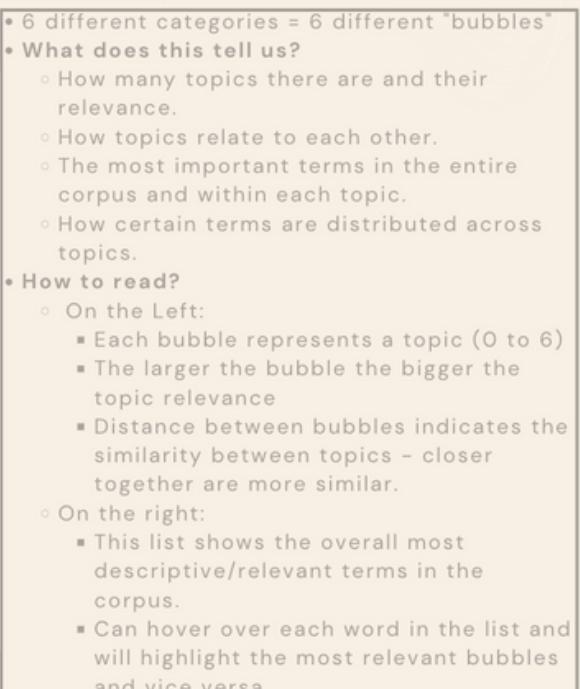


QUESTIONS?

VISUALIZATION PROCESS



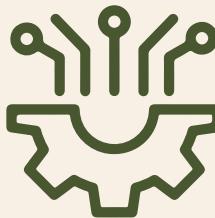
VISUALIZATION PROCESS



SO WHICH METHOD IS BETTER?



Manual



Supervised



Unsupervised

PROS

- Data Familiarization
- Quality Control
- Precision and Accuracy

CONS

- Time-Consuming & Costly
- Oversight/ Inconsistencies
- Human Error: Typos

- High Accuracy
 - Assumption: The Given Data is Also Accurate
 - Strong Predictive Power
 - Data and Time Efficiency

- Time-Consuming & Costly Process
- Oversight/ Inconsistencies:
- Data Processing:

- No Need for Labeled Data
- Great for Data Exploration
- Data and Time Efficiency

- No clear objective with the data
- Sensitivity:
 - Must Filter Through the Words to Avoid
- Data Processing
 - Lots of Time is Needed to Review the Results and Ensure Accuracy