

# Bridging the gap between supervised classification & unsupervised topic modelling for social-media assisted crisis management

## Big Picture: Finetuned Tweet Embeddings (FTE)

**Motivation:** While crisis tweet classification predict useful labels (that may not all generalize to new crises), topic models discover event specific topics (that may not all be useful). Can we bridge this gap?

**Method:** Cluster BERT embeddings learned from supervised tweet classification (FTE) into topics; then extract keywords using tf-idf and attention; measure quality of clusters/keywords with human eval.

**Results:** In a novel snow storm crisis event, relevant classes from the supervised training are preserved, and novel event-specific topics are discovered in an unsupervised way; automatic/human evaluation show that FTE improves over topic models and vanilla BERT embeds.

**Takeaways:** Clustering representations learned from supervised classification can help adapt to domains with overlapping and different latent classes. Human involvement is key to ensure model outputs are aligned with stakeholder needs, especially in crisis management.

## Interpretable Topics for Crisis Managers

Based on annotator interpretations of extracted keywords, FTE clustering preserved relevant classes from supervised training, and discovered novel event-specific topics.

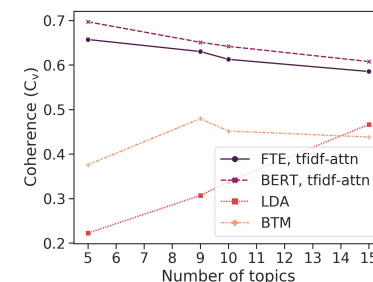
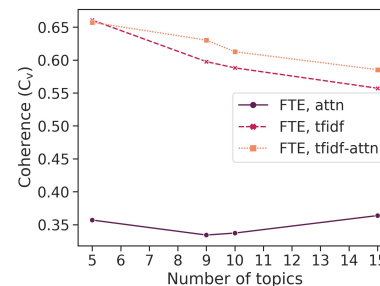
Model	1	2	3	4	Topic 5	6	7	8	9
FTE	reporting monster snowiest recorded peak temperature cloudy reported equivalent meteorologist	ivyparkxadidas mood song blackswan le snowdoor perspective ode music adidasxivypark	outage campus advisory reported impassable remaining thousand requesting suspended	assistance assist troop volunteer providing relief aid request offering rescue	prayer praying pray wish wishing humanity aid surviving loved kindness	blowingsnow drifting advisory caution advised stormsurge wreckhouse surge drifting avoid	trapped stranded hydrant ambulance dead garbage rescue permitted body helped	monster meteorologist drifting perspective stormofcentury mood snowdrift climate windy snowdoor	bread song coffee milk feelin pin enjoying laugh favorite girl
	Weather related	Unrelated information	Power outages	Donation + volunteer	Sympathy + support	Caution + advice	Trapped people		

## Evaluation

Traditional evaluation methods are useful in selecting hyperparameters. However, human evaluation is essential to ensure a model is valuable to crisis managers.

## Automatic Evaluation

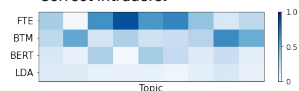
Combining tf-idf and attention for keyword extraction improved coherence in FTE. FTE performed better than traditional topic models, similar to vanilla BERT embeds.



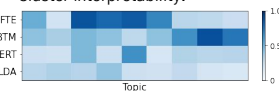
## Human Evaluation

With high agreement, annotators found FTE produces more interpretable, coherent, and useful topics in the context of crisis management. Results for vanilla BERT embeds show limitations of automatic evaluation and need for human-centered ML.

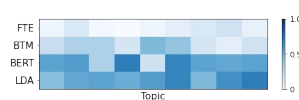
Correct intruders.



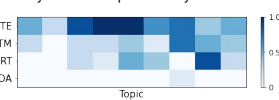
Cluster interpretability.



Unsure intruders.



Keyword interpretability.



Keyword Evaluation scores averaged across topics, number of topics with average scores greater than 0.5.

Score	Average Score	Topic Count	Fleiss' κ			
	BTM	FTE	BTM	FTE		
Interpretability	31.94	<b>65.28</b>	1	<b>5</b>	15.01	<b>17.97</b>
Usefulness	27.78	<b>59.72</b>	1	<b>5</b>	12.36	<b>21.55</b>

Cluster Evaluation scores averaged across top-ics, number of topics with average scores greater than 0.5

Score	Average Score	Topic Count	Fleiss' κ			
	BTM	FTE	BTM	FTE		
Interpretability	50.28	<b>51.53</b>	3	<b>4</b>	11.05	<b>23.45</b>
Usefulness	45.46	<b>46.11</b>	3	<b>5</b>	<b>21.82</b>	21.60
Correct Intruders	35.28	<b>44.17</b>	2	<b>4</b>	25.78	<b>31.50</b>
Unknown Intruders	26.39	<b>8.89</b>	0	0	-	-