Making Text Analytics More Approachable: From Traditional to Generative AI for Business Applications Paper # AL96

Vikranth Reddy Chapaala is a graduate student in the Business Analytics and Data Science program at Oklahoma State University, Stillwater, Oklahoma, with an expected graduation in May 2025. He holds a bachelor's degree in Computer Science and developed a passion for solving meaningful problems through data. Between his bachelor's and master's journey, Vikranth worked for approximately three years in Data Engineering and Analytics teams. During the summer, he worked as a Data Science Intern for a leading payroll and human resources company and is continuing in the same role alongside pursuing his master's degree.



Making Text Analytics More Approachable: From Traditional to Generative AI for Business Applications Paper # AL96

Vikranth Reddy Chapaala, Graduate Student Oklahoma State University Stillwater, Oklahoma



Background





Focus Area

To conduct a comparative study on performance of unsupervised topic modeling techniques:

- Spanning from traditional to generative models
- Devise a new approach for combining topic modeling and text classification.

Techniques Used

- SAS Visual Text Analytics
- Text Analytics with Python



Image Credit: DALL-E



Topic Modeling

Key Idea: Documents are mixtures of latent topics, where a topic is a probability distribution over a word

- The main way of automatically capturing the meaning of documents.
- The topic of an image: a cat, a dog..

Politics, President, Law, Policy.... Space, Planet, Astronaut, Mission.... Sports, Team, Player, Coach, Stadium....







- 2,225 articles published online
- Each article is labeled under one of 5 categories:
 - o Business
 - Entertainment
 - Politics
 - o Sport
 - o Tech

Sample Data		
text	label_text	
wales want rugby league training	Sport	
new harry potter tops book	Entertainment	

https://huggingface.co/datasets/SetFit/bbc-news

http://mlg.ucd.ie/datasets/bbc.html



Topic Modeling & Classification



Evaluations

Topic Modeling

Text Classification

Technique	# Topics	Coherence score
Lda	14	0.34
Top2vec	10	0.47
SAS	11	0.49
Bertopic	53	0.62

Model Type	Accuracy
Transformer based	88.3%
Non- Transformer based	72.6%



Why did the numbers vary?

Traditional

- Based on probabilistic models or matrix factorization.
- Uses bag-of-words assumption (ignores word order).
- Focuses on identifying word cooccurrence patterns in documents.

Modern

- Leverages word embeddings (BERT, GPT) to capture context.
- Uses deep learning and transformers for topic inference.
- Focuses on generating Context-aware topics using semantic relationships.







Conclusions



Transformer-based algorithms improve topic modeling by producing organized clusters and revealing key relationships. Combining topic modeling with language models enhances automatic categorization.



Traditional algorithms often miss text semantics, but large language models (LLMs) provide significant improvements.



Future work will focus on optimizing parameters, conducting detailed studies on specific user groups, and developing scalable platforms.



Thank You!

Vikranth Reddy Chapaala Graduate Student, Business Analytics and Data Science Oklahoma State University vikranth.reddy@okstate.edu (405) 269-2368 https://www.linkedin.com/in/vikranth-reddyc/



Trademark Citation

SAS and all other SAS Institute Inc. product or service names are registered trademarks or trademarks of SAS Institute Inc. in the USA and other countries.
(R) indicates USA registration.

Other brand and product names are trademarks of their respective companies.

