NLP and machine learning to automate identification of suspected medication errors from real-world unstructured narratives

<u>Jeffery L Painter,¹ François Haguinet,² Carolyn Cranfield,³ Andrew Bate³</u> ¹GSK, Durham, NC, USA; ²GSK, Wavre, Belgium; ³GSK, London, UK

Background and aims

Natural language processing (NLP) allows for the transformation numeric feature vectors which can be used in machine learnin



We determined whether or not unstructured te (e.g., safety reports) contain mention of a medi (MedError) by using NLP and ML.

In case a MedError was present, we aimed to su

- error with stated adverse drug reaction (ADR)
- error without harm
- intercepted error
- potential error

This classification is usually rules-based and requires manual review by trained safety scien

Main aim: to build ML models for automating this process, firs narratives contain MedErrors or not, then identifying the corre

Methods

Safety reports were extracted from internal GSK safety databa between January–November 2021 for training our models.



RW, real-world; **TF-IDF**, term frequency-inverse document frequency; **avg**, average; **SQL**, structured query language; **EMA**, European Medicines Agency.

ISPOR 2023, 7–10 May 2023, Boston, MA, USA

	Results			
tion of text into ng (ML) applications. ¹			Binary cla	ssifier p
ext or narratives			Precision	Recal
allation en or	Yes		0.80	0.92
sub-classify it as R)	No		0.88	0.73
	Accuracy			
	Macro avg		0.84	0.82
	Weighted a	avg	0.84	0.83
ntists. est identifying if rect sub-classification.	A A A Potential Vithout harm			
	Intercepted			
	With ADR			
ase				
ta cleaning rratives were processed ng text vectorization ² -IDF ²), and all product		A Beri (bina) Split of a 75%	noulli naïve ry) classifie data (RW unst	e Bayes er ^{&} was t ructured (25%



Steps to building an NLP and ML model for MedErrors identification

STEP 1	STEP 2
 Data acquisition and formatting case narrative data for use in model building Data extraction from the ARGUS Oracle database using SQL to develop training and test data sets. 	 Data labeling Narratives were labeled according to their binary class: case or control (not a MedError). For a case, a secondary label indicating the sub-class of MedError found (e.g., with ADR, without harm) was applied.

To minimize bias in labeling, a random sample of 3,122 case narratives (containing MedErrors) for ML training and 1,611 controls was collected. ^{}Chosen as the basis of our model building due to its general performance characteristics⁴.

Funding: GlaxoSmithKline Biologicals SA. Acknowledgments: Medical writing (Lucia Adina Truta), design and coordination support were provided by Akkodis Belgium c/o GSK. Disclosures, References available via the QR code.





performance

STEP 3

Data cleaning

NLP methods generally perform better after "cleaning" the data. A first test showed that product names weighted too much in the multi-label classification. Therefore, the product names were masked for the model building.



Performance was evaluated using cross-validation.



Concordance of labeling (based on a rules-based approach) was compared to an unsupervised clustering method using k-means for narratives not sub-classified.

STEP 4

- Building a binary classifier
- (1) Narrative represents a MedError or
- (2) Narrative does not represent a MedError
- Building a multi-class identifier classifier (conditioned to the

STEP 5

binary classifier [step 4]) If a narrative was identified as MedError

in step 4, next step was to attempt further sub-classification into the EMA categories (e.g., with ADR, without harm).