UTILIZING A TRANSPARENCY-DRIVEN ENVIRONMENT TOWARD TRUSTED AUTOMATIC GENRE CLASSIFICATION

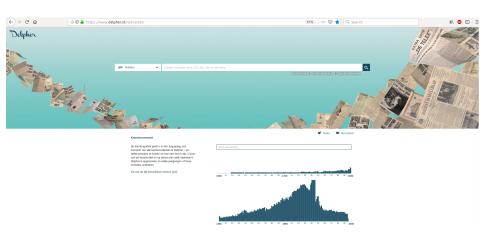
A Case Study in Journalism History

Aysenur Bilgin ¹, **Erik Tjong Kim Sang**², Kim Smeenk³, Laura Hollink¹, Jacco van Ossenbrugen¹, Frank Harbers³, Marcel Broersma³ November 1, 2018

¹CWI, ²Netherlands eScience Center, ³University of Groningen

aysenur.bilgin@cwi.nl

NEWSPAPER ARCHIVES



Newspaper archive Delpher.nl offers millions of newspaper pages

MOTIVATION

- · Study genre distribution in Dutch newspaper text
- · However, digital archives do not contain genre labelling
- · Due to big corpus size, we need machine learning

BUT... Can we trust the results of machine learning?

TRANSPARENCY

- \cdot We need to be able to show what machine learning is doing
- · Both for specialists and non-specialists (domain experts)
- · We want transparency in all aspects of the labelling process

TRANSPARENT WORKFLOW

Data pre-processing

Upload formatted data, choose representation type and pre-processing steps.

Individual pipeline analysis

Using several metrics, receive insights about the performance of the pipeline (e.g. accuracy, precision, etc.) and its global interpretations if available.

2 Optimal model selection

Optional step to detect the optimal machine learning algorithm and its hyperparameters.

Pipeline comparison

Based on test data, compare collective predictions of the pipelines in tabular views and get insights about the pipelines using local explanations.

3 Configuration and training

Use the optimal choices from the previous step if applied or configure the algorithm and train.

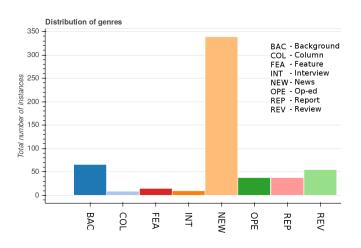
6 Domain hypothesis testing

Choose the best suited pipeline for the research hypothesis and apply it to large-scale real-world data.

CASE STUDY

- We study distribution of genre (e.g. column, interview, review, etc.)
 over time
- · We use supervised machine learning algorithms for genre labelling
- · We can analyze the impact of different variants for:
 - · Data preparation
 - · Data pre-processing
 - · Machine learning algorithms
 - · Algorithm configurations
- \cdot We examine visualizations for fair comparison

DATA PREPARATION



We need a balanced dataset to ensure a fair model

DATA PRE-PROCESSING

We consider two types of document representations:

· Bag-of-words (i.e. words and their frequencies)

· Linguistic features (e.g. document length, number of pronouns, etc.)

MACHINE LEARNING ALGORITHMS AND CONFIGURATIONS

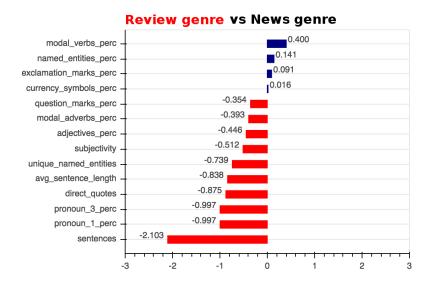
We evaluate three machine learning algorithms:

- · Naive Bayes
- · Support Vector Classifier
- · Random Forest

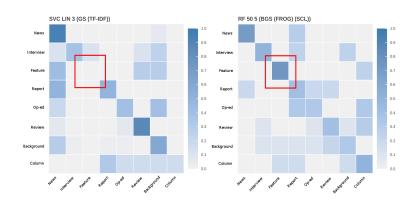
The algorithms have several configurable parameters that have impact on their behaviour.

We use automated techniques to select optimal values of these parameters.

PIPELINE ANALYSIS - FEATURE IMPORTANCE RANKING



PIPELINE COMPARISON USING CONFUSION MATRICES



Confusion matrices for two machine learning pipelines: SVC accuracy is 70% (left) and RF accuracy is 41% (right)

PIPELINE COMPARISON USING ARTICLE-BASED VIEW

Article #	Article text	Mutually Agreeing Pipelines 👢	Prevailing Genre (prediction)	True Genre ↓↑
1800	C.P.N . jaren van geharrewar is het er he ff 6 $^{\circ}$ S $^{\circ}$ tt : e van gekomen	RF 1000 5 (GS (FROG) (SCL)) RF 50 5 (BGS (FROG) (SCL)) RF 1000 2 (CGS (FROG) (SCL)) SVC LIN 3 (CGS (TF-IDF) (SWR))	News	Op-ed
55	Produktschap gaat praten over vleesprijzen (Van onze Haagse redacteur)	RF 1000 5 (GS (FROG) (SCL)) RF 50 5 (BGS (FROG) (SCL)) RF 1000 2 (CGS (FROG) (SCL)) SVC LIN 3 (CGS (TF-IDF) (SWR)) SVC LIN 3 (CGS (TF-IDF))	News	News
77	Artis schenkt Blijdorp orang oetang De directie van diergaarde Blijdorp	RF 1000 5 (GS (FROG) (SCL)) RF 50 5 (BGS (FROG) (SCL)) RF 1000 2 (CGS (FROG) (SCL)) SVC LIN 3 (CGS (TF-IDF) (SWR)) SVC LIN 3 (CGS (TF-IDF))	News	News
111	Ingestudeerde spelpatronen komen niet uit de verf tegen Oostenrijk Door	RF 1000 5 (GS (FROG) (SCL)) RF 50 5 (BGS (FROG) (SCL)) RF 1000 2 (CGS (FROG) (SCL)) SVC LIN 3 (CGS (TF-IDF) (SWR)) SVC LIN 3 (CGS (TF-IDF))	Background	Background
117	Simpele zege volleyballers op Portugal SOLBJERG , 2 april — Het	RF 1000 5 (GS (FROG) (SCL)) RF 50 5 (BGS (FROG) (SCL)) RF 1000 2 (CGS (FROG) (SCL)) SVC LIN 3 (CGS (TF-IDF) (SWR)) SVC LIN 3 (CGS (TF-IDF))	News	News

PIPELINE COMPARISON USING SET-BASED VIEW

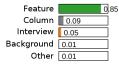
Mutually Agreeing Pipelines 🕸	Predicted Genre 🚉	Number of Articles ↓ ₽	Articles
SVC LIN 2 (BGS (TF-IDF) (SWRI)) SVC RBF 10 (BGS (TF-IDF))	Report	208	[Article 13] [Article 73] [Article 28] [Article 49] [Article 49] [Article 58] [Article 58] [Article 58] [Article 175] [Article 1
RF 1000 5 (GS (FROG) (SCL)) RF 1000 2 (CSG (FROG) (SCL)) RF 50 5 (BSS (FROG) (SCL)) SVC LIN 3 (CSG (TF-IDF)) SVC LIN 3 (CSG (TF-IDF) SVC LIN 3 (CSG (TF-IDF)) SVC LIN 3 (CSG (TF-IDF))	News	142	Article 49 Article 49 Article 49 Article 51 Article 150 Article 371 Article 272 Article 278 Article 286 Article 286 Article 286 Article 286 Article 286 Article 286 Article 371 Article 372 Article 375 Article 397 Article

PIPELINE COMPARISON USING EXPLANATION-BASED VIEW

Article # ↓≟	Pipeline 11	Predicted Genre ↓↑	Text Explanation ↓↑	Feature Explanation 11
1	RF 1000 2 (CGS (FROG) (SCL))	Feature	N/A	Click here for feature explanation
59	RF 1000 5 (GS (FROG) (SCL))	Feature	N/A	Click here for feature explanation
59	RF 1000 2 (CGS (FROG) (SCL))	Feature	N/A	Click here for feature explanation
59	RF 50 5 (BGS (FROG) (SCL))	Feature	N/A	Click here for feature explanation
66	RF 50 5 (BGS (FROG) (SCL))	Feature	N/A	Click here for feature explanation
80	RF 1000 5 (GS (FROG) (SCL))	Feature	N/A	Click here for feature explanation
80	RF 1000 2 (CGS (FROG) (SCL))	Feature	N/A	Click here for feature explanation
87	RF 1000 5 (GS (FROG) (SCL))	Feature	N/A	Click here for feature explanation
94	RF 1000 5 (GS (FROG) (SCL))	Feature	N/A	Click here for feature explanation

LOCAL EXPLANATIONS

Prediction probabilities



NOT Feature Feature



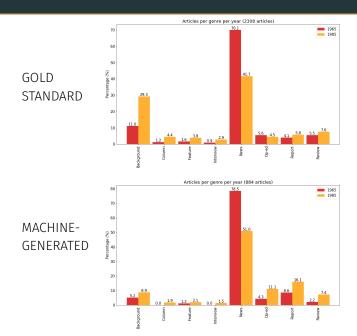
Prediction probabilities



NOT Feature Feature

reattile
tokens > 1053.00
li.04
sentences > 55.25
li.04
direct_quotes > 12.00
lo.02
pronoun_2_perc > 0.00
lo.02
adjectives > 82.00
lo.02
pronoun_3 > 22.00
lo.02
cogn_verbs_perc <= 0.00
lo.01
pronoun_2 > 3.00
lo.01
lo.01

DOMAIN HYPOTHESIS: GENRE DISTRIBUTION OVER TIME



CONCLUDING REMARKS

- · With accuracy alone, we cannot select the best pipeline for the domain scientist
- · Transparency promises in-depth insights and discussions for both the domain scientist and the computer scientist
- · We can claim transparency increases trust of the domain scientist in the application of machine learning

UTILIZING A TRANSPARENCY-DRIVEN ENVIRONMENT TOWARD TRUSTED AUTOMATIC GENRE CLASSIFICATION

A Case Study in Journalism History

Aysenur Bilgin ¹, **Erik Tjong Kim Sang**², Kim Smeenk³, Laura Hollink¹, Jacco van Ossenbrugen¹, Frank Harbers³, Marcel Broersma³ November 1, 2018

¹CWI, ²Netherlands eScience Center, ³University of Groningen

aysenur.bilgin@cwi.nl