

# Evaluating UK Honours Candidates using a Novel Data-Analytics Pipeline

Francesca von Braun-Bates<sup>1\*</sup>, Sunreeta Sen<sup>2</sup>, Indraayudh Talukdar<sup>3</sup>,  
Anirban Lahiri<sup>4</sup>

<sup>1\*</sup>Ministry of Justice, Government of United Kingdom, London, United Kingdom.

<sup>2</sup>Arndit Ltd., Cambridge, United Kingdom.

<sup>3</sup>Indian Institute of Technology Delhi, New Delhi, India.

<sup>4</sup>Kainos, London, United Kingdom.

## Abstract

This paper explores the application of data science to the UK Honours system, focusing on how publicly available information can support due diligence and enhance public confidence. We identify key challenges that must be addressed before such methods can be applied responsibly in this context. We present a data-processing pipeline that collects open-source information via web scraping, extracts relevant text using coreference resolution, and applies sentiment analysis to identify potentially salient evidence about candidates. We evaluate two established sentiment analysis algorithms (AFINN, VADER) and introduce a novel method, MINOS, tailored to this use case.

Our results indicate that sentence-level filtering and domain-specific sentiment modelling improve the identification of relevant positive and negative signals. The proposed system is intended to augment, rather than replace, human judgement by highlighting potentially high-impact information for further review. This approach may help reduce the risk of oversight in candidate assessment and can be extended to other awards and recognition processes.

**Keywords:** Sentiment Analysis, Natural Language Processing, Open-Source Intelligence, Decision Support Systems

## 1 Introduction

It is always an honour for a person to be nominated for an award. Since 1066, the monarch of the United Kingdom has recognised pre-eminent citizens by awarding Honours [1]. These stand out amongst national awards (see Annex 2 of [2] for a comparison) for their longevity, adaptability to keep the system relevant and accessible in the modern age, and the wide range of achievements rewarded. Since these awards are highly coveted, thorough due diligence needs to be done on each applicant especially the prospective candidates. Failing this, the general public might lose faith in the fairness of the Honours system. So far this entire process of due diligence would be done manually by officials. Furthermore, a core principle of the Honours System is that the award must be maintained via continued good conduct by the recipient [3]. Should the recipient's behaviour "bring the Honours system into disrepute", they may be stripped of their Honour to preserve public confidence in the system, a process known as "forfeiture" [3]. The obvious question, then, is how to do this fairly, robustly, and without undue public expense? With over 150 000 recipients of the Order of the British Empire alone [2], this is an arduous task and can benefit from leveraging the power of modern data science techniques. We examine the challenges in this unique problem in this paper. At the same time, this work tries to answer the questions:

*How to measure public sentiment about a living individual?*

*How to define and identify "red flag" behaviours with open-source intelligence?*

This would in turn have wide-ranging applications including award decisions for many other honours and awards, corporate Human-resource management, investigative journalism to national security vetting.

To reduce human intervention and produce impartiality, the procedure needs to be automated. This paper does not suggest *replacing* human selection. Instead our methods would augment the current human decision-making by extracting relevant information from data in an auditable and robust way. The extraction of insight would remain the human committees' responsibility.

The goals for this work are:

- **Maximise available data:** Use web-scraping to scan more of the web than a human could do, more efficiently, auditably and responsibly.
- **Minimise human sifting:** Use algorithms to identify relevant information from our large volume of data.
- **Identify key evidence:** Both the most positive and negative assertions can be flagged by looking at extremes of sentiment.
- **Compare multiple viewpoints:** Provide multiple sentiment algorithms so that assessment does not have a single point of failure.
- **Ensure attributability:** Ensure that both the original article content and the source URL are preserved in our results. This enables a human sifter to examine the source, its biases and quality, to judge any (positive or negative) evidence.

Although this appears to be a classic data science application—to extract a small amount of relevant information from a large corpus—we need careful handling to avoid misleading results. Three different sentiment analysis algorithms have been applied to web-scraping of search results for a given subject of interest. Our subjects of interest (hereafter referred to as SOIs) are drawn from three distinct categories – Infamous individuals also referred to as test cases, successful Honours recipients, and Honours forfeiture cases – to illustrate the complexity in classifying *algorithmically* what to a human are three highly distinct behavioural groups.

Some of the key challenges addressed in this work:

- Search engines do not necessarily rank their results in order of relevance to their search query, which biases the impressions of a human assessor.
- Judging relevance at an article level leads to misleading sentiment results; only at the individual sentence level can we guarantee relevance.
- Evaluation of Co-reference resolution algorithms to extract all sentiment-bearing sentences about a given individual, rather than considering irrelevant information.
- Off-the-shelf algorithms have difficulty identifying “red flag” behaviours which would lead to forfeiture, even in the test group.
- Improving on the quality of the result — rather than increasing the complexity of the existing algorithms — therefore a simple approach was developed from first principles.

Hence, the three core components of the proposed method are web-scraping, coreference resolution and sentiment analysis.

Web-scraping has existed in various forms since the earliest days of the internet [4] including web crawling, web indexing and archiving of web content. We use “web-scraping” to refer specifically to the extraction of raw text from the (non-dark) web. Once extracted, the text is cleaned, chunked, and piped to a sentiment analyser. This process is explained in Section ??.

Coreference resolution is the task of associating different linguistic expressions which refer to the same entity [5]. It is a fundamental building block for more complex NLP tasks. We applied this to trace references to an SOI throughout an article. In this way we identified sentences which specifically involved our individuals, examining only the sentiment from those sentences to minimise noise in the results.

Many sentiment analysis approaches exist, and it is well beyond the scope of this paper to canvas them (see [6] for one such attempt). We apply the popular VADER model [7] optimised for social media sentiment analysis. We also use the simpler AFINN model, which is a composite of existing lists and some internet slang [8]. We define a custom sentiment model MINOS which we benchmark against two operating modes of VADER and the results of AFINN. We summarise all three algorithms in ??.

Sentiment analysis of public figures appears to be (curiosity) limited both in the scope of source data and the depth of analysis. This paper extends the *status quo* in three ways:

- it uses a broad range of sources rather than purely social media or news articles;
- a custom sentiment algorithm has been utilised alongside using two existing methods [7, 8].
- the relevance of any text input is validated using coreference resolution.

- SOIs from many different fields and with varying degrees of celebrity (or non-celebrity) status have been analysed.

The research in [9] uses the fraction of positive tweets by 400 celebrities as a proxy for the individuals’ positivity or negativity. This differs from our method, where we measure their sentiment based upon open-source information *about* the public figures themselves, not necessarily (nor solely) produced *by* them. The paper [10] applies TextBlob’s polarity measures to tweets returned from the query “Elon Musk”, as an estimate of how Twitter users felt about Musk’s proposed acquisition of the Twitter platform. The only academic paper we found similar to this work is [11], which sources large-scale European news data and estimates the sentiment of named entities (including public figures). However, this paper [11] aims to measure reporting bias amongst news sources—by assuming that their sources are biased, they preclude themselves from using news stories to obtain an unbiased measure of public sentiment. We make the opposite assumption, i.e. that our sources represent genuine differences in opinion or report different subsets of (truthful) facts.

The structure of this paper is as follows. In Section ?? we summarise the aspects of the Honours system which are relevant to our forfeiture problem and define our SOIs. Our methodology is described in Section ??, including data collection, cleaning, sentiment analysis and extraction of results. Section ?? shows the results thus highlighting the advantage of web-scraping over manual research by examining the fraction of web-scraped text which survived our relevance filter criteria. Thereafter the results corresponding to various options and alternatives for the algorithm using coreference resolution have been illustrated. This is followed by the final results which can be used by analysts to determine if a candidate is suitable for an award or not. The finally concludes (5) with the contributions of this work and potential future work.

## 2 Background

This section provides minimal domain context and defines the subjects of interest (SOIs) used to evaluate our methodology.

### 2.1 The Honours system and forfeiture

The UK Honours system recognises individuals for exceptional contributions across a range of fields. Awards are based on evidence of sustained achievement, public impact, and service beyond an individual’s primary role. Once awarded, an Honour may be revoked if subsequent conduct is deemed to bring the system into disrepute, a process known as forfeiture.

From a data analysis perspective, this creates a classification challenge: individuals may exhibit predominantly positive, predominantly negative, or mixed public sentiment over time. In particular, forfeiture cases represent individuals whose publicly observable behaviour changed significantly after recognition.

### 2.2 Forfeiture as a data analysis problem

Forfeiture decisions are not governed by fixed rules or deterministic thresholds. Instead, they are made on a case-by-case basis, based on qualitative judgements about whether an individual’s conduct brings the Honours system into disrepute [3, 12]. For example, forfeiture is almost certain in the following cases [3, 13, 14]:

- disbarment from a professional body in the field in which the recipient was nominated
- censure by a regulator if directly relevant to the nomination
- criminal conviction resulting in a prison sentence of three months or more
- any criminal offence under the Sexual Offences Act 2003 and related legislation

However, it is not guaranteed, and past cases do not create an automatic precedent. This lack of a well-defined decision boundary makes the problem difficult to formalise as a conventional classification task.

In addition, observed forfeiture outcomes are sparse and often occur long after the underlying behaviour. Comparing the annual number of awards to forfeitures in *e.g.* Annex B of [15] show that this is a sub-percent occurrence. Table A3 shows the lag between award and forfeiture for our SOIs can be decades. Relevant evidence is expressed across heterogeneous sources and in indirect forms, making it difficult to identify using document-level or keyword-based approaches.

The problem is further complicated by the temporal evolution of an individual’s public profile. Individuals subject to forfeiture frequently exhibit both positive and negative sentiment over time, leading to

mixed or non-stationary signal distributions. Finally, the asymmetric cost of errors—where missing negative evidence is more consequential than identifying spurious signals—motivates a conservative approach that prioritises the detection of potentially adverse information.

## 2.3 Subjects of interest

To evaluate the proposed methodology, we construct three groups of SOIs:

- **Awarded:** individuals who have received an Honour and retained it, representing consistently positive public records.
- **Infamous:** individuals associated with serious misconduct or criminal activity, representing strongly negative cases.
- **Forfeited:** individuals who were awarded an Honour and later had it revoked, representing mixed or evolving public sentiment.

Each group contains twenty individuals drawn from a range of domains and levels of public prominence. These groups are intended as a structured test set rather than a representative sample of the population.

The “infamous” group provides a negative extreme, where strong negative sentiment is expected. The “awarded” group provides a positive extreme, where sentiment should be consistently favourable. The “forfeited” group represents the primary use case, in which both positive and negative signals are present. A successful methodology should distinguish between these groups and identify cases where negative evidence is sufficiently prominent to warrant further investigation.

This framing allows us to assess whether sentiment signals extracted from open-source data can provide meaningful support for identifying high-risk or anomalous cases within a broader population.

## 3 Methodology

This section describes a data-processing pipeline designed to extract sentiment signals about a subject of interest (SOI) from heterogeneous open-source text. The pipeline comprises: (i) data collection, (ii) relevance filtering, (iii) sentence-level extraction with coreference resolution, (iv) sentiment analysis using baseline and domain-specific models, and (v) aggregation of results.

### 3.1 Data collection and processing

For each SOI, we query a search engine using the individual’s full name and collect the top ~60–70 results. From each result, we extract raw text and retain the source URL to enable traceability and subsequent human verification [16, 17].

We restrict analysis to publicly accessible content on the open web and do not access paywalled or authenticated sources. This introduces sampling bias but reflects realistic constraints faced by human evaluators and common limitations of web-scraping approaches [4, 18, 19].

All text is tokenised at the sentence level using the NLTK toolkit, as sentiment is evaluated at the level of individual statements rather than entire documents [20, 21].

### 3.2 Article relevance criteria

The extracted text contains a substantial proportion of irrelevant or weakly related content. We therefore apply two simple filters as proxies for relevance:

1. File size  $\geq$  500 characters;
2. The SOI’s name appears at least once.

The file size threshold removes trivial or non-content pages (e.g. navigation or login prompts), while the name filter ensures a minimal level of topical relevance. Although naming conventions introduce edge cases (e.g. titles, multiple surnames, or alternative forms), a simple heuristic based on commonly used names was found to be sufficient for filtering purposes.

### 3.3 Coreference resolution

Document-level filtering is insufficient to guarantee that sentiment-bearing text refers to the SOI. Even within relevant documents, many sentences describe other entities or general context. To address this, we

apply coreference resolution to identify expressions that refer to the SOI and restrict analysis to those sentences [5, 22].

Coreference resolution replaces pronouns and other referring expressions with their corresponding named entities, allowing indirect references to the SOI to be captured. This reduces false negatives (relevant sentences excluded) and improves the precision of sentiment extraction, though the process is inherently probabilistic and may introduce minor errors.

### 3.4 Sentiment analysis

We evaluate two established sentiment analysis methods alongside a domain-specific model.

AFINN is a lexicon-based approach that assigns integer sentiment scores to individual words [8, 23, 24]. VADER extends this approach with rules that account for negation, emphasis, and punctuation, producing both component and compound sentiment scores [7, 25]. Both methods are designed primarily for short, informal text and can perform poorly on heterogeneous, long-form sources.

To address this limitation, we introduce MINOS, a domain-specific sentiment model designed to prioritise the detection of negative signals. MINOS uses separate positive and negative lexicons derived from multiple sources, including domain-relevant terminology (e.g. criminal activity and honours-related terms). Words not present in either lexicon are treated as neutral.

The key design principle of MINOS is that negative evidence dominates: if a sentence contains any negative term, the overall sentence sentiment is treated as negative. This reflects the asymmetric cost of errors in the application context, where missing negative signals is more consequential than identifying spurious positives.

### 3.5 Aggregation of sentiment

Each sentence is assigned a sentiment score based on the selected model. We summarise sentiment for an SOI by computing the arithmetic mean of non-zero sentence scores, excluding neutral statements that do not contribute meaningful signal.

Sentences are treated as the fundamental unit of analysis—rather than words, paragraphs or articles. For a given SOI, the aggregate sentiment is therefore obtained by integrating over the distribution of sentence-level scores across all “relevant text” (where the level of relevance is controlled by the filtering). The resulting value represents the centroid of this distribution.

This sentence-level aggregation is necessary because article-level scoring is not equivalent across models. For lexicon-based approaches such as AFINN and VADER, sentiment can be accumulated across words or sentences. However, for MINOS, which assigns negative sentiment whenever any negative term is present, aggregating at the article level would collapse mixed-sentiment articles into purely negative scores. This would obscure the presence of both positive and negative evidence within the same document.

The full distribution of sentence-level scores therefore provides additional information beyond the centroid, particularly in cases where sentiment is mixed. In the results, we analyse both the centroid and the distribution to characterise differences between SOI groups.

The sentiment models considered in this work operate on different scales and scoring schemes. As a result, their outputs are not directly comparable in absolute terms. Instead, comparisons are made within each model by examining the relative separation between SOI groups and the structure of their sentiment distributions. This allows us to assess how effectively each method distinguishes between positive, negative, and mixed sentiment profiles without requiring cross-model normalisation.

### 3.6 Output for human analysis

The objective of the pipeline is to support human evaluation by identifying potentially relevant positive and negative evidence within large text corpora. In addition to aggregate sentiment measures, the system retains the underlying sentences and their source URLs, allowing analysts to inspect high-impact evidence directly and assess its reliability and context.

## 4 Results, Analysis and Discussion

### 4.1 Overview

This section evaluates the proposed methodology with respect to the challenges identified in Section 2. In particular, we assess the extent to which the pipeline is able to extract meaningful sentiment signals from noisy, heterogeneous open-source data and distinguish between different categories of subjects of interest (SOIs).

Recall that analysis is structured around three SOI groups: awarded, infamous, and forfeited. These groups represent positive, negative, and mixed sentiment profiles respectively, providing a basis for evaluating how well the methodology captures different types of behavioural signals.

The SOI groupings provide proxy-labelled test cases for evaluation; the objective is not to train a model to classify individuals, but to assess whether the extracted sentiment signals are consistent with known characteristics of each group.

We focus on three key aspects of performance. First, we examine the impact of relevance filtering on reducing noise in the input data. Second, we compare the behaviour of baseline sentiment models with the proposed MINOS approach. Finally, we analyse the resulting sentiment distributions to assess how effectively the method distinguishes between SOI groups.

## 4.2 Addressing signal heterogeneity: relevance filtering

A primary challenge identified in Section 2 is the heterogeneous and weakly relevant nature of open-source text. Web search results contain large volumes of content that are either unrelated to the SOI or contain sentiment that does not refer to the individual. This subsection evaluates the effectiveness of the proposed filtering steps in reducing this noise.

**fig:filesize-distribution (placeholder)** shows the distribution of file sizes for retrieved articles across all SOIs. The results indicate substantial variation in article length, spanning several orders of magnitude. A significant proportion of files are either too short to contain meaningful information (e.g. navigation pages or login prompts) or sufficiently large that they are unlikely to be focused on a single individual.

Applying a file size threshold removes a large number of trivial or spurious results. However, file size alone is insufficient to ensure relevance. **fig:filter-survival (placeholder)** shows the proportion of articles that pass the file size and name-based filters. In all SOI groups, a substantial fraction of retrieved content fails at least one of these criteria, demonstrating the extent of noise in the raw data.

Name-based filtering further reduces irrelevant content by requiring explicit mention of the SOI. However, even after these filters are applied, many retained articles contain sentences that are not directly related to the individual. This motivates the use of sentence-level filtering and coreference resolution, which are evaluated in subsequent sections.

Overall, these results demonstrate that relevance filtering is a necessary preprocessing step for extracting meaningful sentiment signals. Without such filtering, sentiment analysis would be dominated by noise arising from unrelated or weakly relevant text.

## 4.3 Baseline model limitations

Weak signals in large, unstructured data are only partially enhanced by the filtering mechanisms in the previous **sub:filtering-relevance**. Sentiment analysis models must still distinguish between sentiment that genuinely refers to the SOI and sentiment arising from unrelated context. This subsection evaluates the behaviour of baseline sentiment models, AFINN and VADER, under these conditions.

We first consider sentiment computed directly from raw, unfiltered data. **fig:centroid-raw (placeholder)** shows the distribution of centroid values for each SOI group using article-level sentiment. The results exhibit substantial overlap between the awarded, infamous, and forfeited groups. In particular, the infamous group, which should be strongly negative, often appears close to neutral. This indicates that sentiment from unrelated or weakly relevant text dominates the aggregate signal.

Applying article-level filtering reduces this effect but does not fully resolve it. **fig:centroid-filtered (placeholder)** shows the corresponding centroid distributions after filtering by file size and name occurrence. While the separation between groups improves slightly, significant overlap remains. This suggests that document-level filtering is insufficient to isolate sentiment that is specifically attributable to the SOI.

We therefore examine sentence-level sentiment, restricting analysis to sentences that explicitly contain the SOI’s name. **fig:centroid-sentence (placeholder)** shows the resulting centroid distributions. This further reduces noise and improves separation between groups. However, even at the sentence level, both AFINN and VADER continue to assign positive sentiment in contexts that are not directly related to the SOI’s behaviour.

This effect arises because baseline sentiment models operate on local lexical cues without resolving which entity a sentiment-bearing expression refers to. While coreference resolution allows us to identify sentences that mention the SOI, it does not determine whether the sentiment expressed in those sentences is attributable to the individual.

As a result, sentiment associated with surrounding context—such as events, topics, or other entities—can be incorrectly attributed to the SOI. For example, Dame Mary Beard’s sentiment scores are

negatively affected by “death” appearing in the title of her work *Pompeii: Life and Death in a Roman Town*, even though it does not describe her character.

This limitation leads to systematic misattribution of sentiment, particularly in heterogeneous text where descriptive or contextual language is common.

Overall, these results demonstrate that standard sentiment analysis methods are not well suited to this task without additional constraints. Although filtering and sentence-level analysis reduce noise, they do not eliminate the influence of irrelevant sentiment. This limitation motivates the need for a domain-specific approach that explicitly accounts for the asymmetric importance of negative signals, which we address in the following subsection.

#### 4.4 Addressing negative recall risk with MINOS

The limitations identified in Section 4.3 arise in part from the inability of baseline sentiment models to attribute sentiment to the correct entity. In addition, the application context introduces asymmetric error costs: failing to identify relevant negative evidence (false negatives) is more consequential than incorrectly flagging spurious signals (false positives). The objective is therefore not to produce an unbiased estimate of sentiment, but to maximise the likelihood of identifying potentially adverse evidence for subsequent human review.

This objective manifests differently across SOI groups. For individuals in the “awarded” group, the goal is to maximise sensitivity to potential forfeiture signals. Even a small number of negative statements may be important, and should be surfaced for triage by a human analyst. Missing such signals represents a failure of the system. Conversely, for individuals in the “infamous” group, we expect overwhelmingly negative sentiment. In this case, the presence of positive or neutral sentiment in the output indicates contamination from irrelevant context, which obscures the true signal.

We first consider a representative case from the “awarded” group. **fig:marks-afinn (placeholder)**, **fig:marks-vader (placeholder)**, and **fig:marks-minos (placeholder)** show the distribution of sentence-level sentiment scores for the same SOI under AFINN, VADER, and MINOS respectively. The baseline models produce a broad distribution with both positive and negative values. While this reflects the presence of diverse language in the corpus, it also indicates that negative signals are diluted by unrelated positive context. As a result, potentially relevant negative evidence is not clearly distinguished from background noise.

We now consider a corresponding example from the “infamous” group. **fig:kasab-afinn (placeholder)**, **fig:kasab-vader (placeholder)**, and **fig:kasab-minos (placeholder)** show the equivalent distributions for this SOI. Despite the strongly negative nature of the underlying behaviour, both AFINN and VADER assign a substantial number of positive or near-neutral scores. This uncovers an axiomatic principle of both AFINN and VADER: that they were designed to estimate overall sentiment without privileging positive or negative evidence. As a result, positive language in surrounding context—such as descriptions of institutions or actions of other individuals—can offset genuinely negative signals.

The contrast between the two cases highlights a fundamental limitation of baseline sentiment models: their design objective is to provide an unbiased estimate of sentiment, rather than to prioritise the detection of specific types of evidence. In this application, such neutrality is undesirable, as it increases the likelihood of false negatives by allowing negative signals to be masked by unrelated positive context.

The MINOS model addresses this limitation by explicitly prioritising negative evidence. For both SOIs, MINOS assigns negative sentiment whenever a sentence contains a negative term, regardless of the presence of positive language. This increases the likelihood that potentially relevant negative statements are surfaced. In the “awarded” example, this results in a small number of clearly identifiable negative signals, suitable for human review. In the “infamous” example, it produces a distribution that is overwhelmingly negative, reducing the influence of irrelevant positive context.

Taken together, these results demonstrate that prioritising the detection of negative evidence provides a more appropriate representation of risk in this setting. While this approach may increase the number of sentences flagged for review, it reduces the likelihood that important signals are missed, aligning the behaviour of the model with the requirements of the application.

#### 4.5 Mixed sentiment profiles and temporal behaviour

A further challenge identified in **sec:forfeiture-challenges (placeholder)** is that an individual’s public profile may evolve over time, resulting in mixed or non-stationary sentiment signals. This is particularly relevant for forfeiture cases, where individuals may have a substantial body of positive reporting prior to the events that led to the withdrawal of an honour. In such cases, sentiment cannot be characterised by a single polarity, and aggregate measures must be interpreted with care.

We therefore examine the distribution of sentence-level sentiment scores for individuals in the “forfeited” group. **fig:vasco-knight (placeholder)** and **fig:rodale (placeholder)** show representative examples under AFINN, VADER, and MINOS. In both cases, the distributions exhibit a mixture of positive and negative scores, reflecting the coexistence of positive reporting associated with earlier achievements and negative reporting associated with subsequent events.

This behaviour contrasts with the more homogeneous distributions observed in the “awarded” and “infamous” groups. For example, **fig:hawking (placeholder)** shows a predominantly positive distribution, while **fig:epstein (placeholder)** shows a predominantly negative distribution. These cases illustrate that, in the absence of temporal shifts in behaviour, sentiment signals tend to be concentrated around a single polarity.

The limitations of centroid-based summaries become particularly apparent when comparing individuals with similar aggregate values. **fig:forfeited-centroid-comparison (placeholder)** shows examples of SOIs with comparable centroid values but markedly different underlying distributions. In such cases, the centroid obscures the presence of both strongly positive and strongly negative evidence, and therefore fails to capture the qualitative differences between individuals.

Examining the full distribution of sentence-level scores provides a more informative representation. In particular, the presence of distinct positive and negative components allows forfeited individuals to be distinguished from both extremes. This supports the use of distributional characteristics, such as spread or multi-modality, in addition to aggregate measures when identifying potentially anomalous cases.

An additional consideration is the cumulative nature of the underlying data. The corpus for each SOI is constructed from publicly available text collected at a single point in time, but reflects information that has accumulated over the individual’s public history. New information is therefore added to an existing body of evidence rather than replacing it. In principle, repeated application of the pipeline as new data becomes available would lead to convergence of both the centroid and the distribution, unless the new information represents a significant deviation from prior behaviour.

In this sense, the methodology captures behaviour integrated over time, partially mitigating the effect of delays in the emergence of negative evidence. While forfeiture decisions may occur after a lag, the underlying signals can be incorporated into the analysis as soon as they appear in publicly available sources. This highlights a potential advantage of automated approaches over manual review, which may be limited in its ability to aggregate historical evidence at scale.

Overall, these results demonstrate that forfeiture cases are characterised by mixed sentiment profiles that cannot be adequately captured by aggregate measures alone. A distributional approach is therefore necessary to identify individuals whose public record contains both positive and negative signals, and to distinguish these cases from those with consistently positive or consistently negative sentiment.

## 4.6 Interpreting outputs under limited and non-deterministic ground truth

A central challenge identified in **sec:forfeiture-challenges (placeholder)** is that forfeiture outcomes are not governed by a fixed or observable decision rule. Instead, they arise from case-by-case judgements based on qualitative assessment of an individual’s conduct. In addition, the number of publicly documented forfeiture cases is limited relative to the total number of honours awarded. Together, these factors mean that there is no well-defined ground-truth label that can be used to train or evaluate a conventional supervised classification model.

The SOI groupings used in this work therefore serve as proxy-labelled test cases for evaluation rather than as training data. The objective is not to assign individuals to these categories with high accuracy, but to assess whether the sentiment signals extracted by the pipeline are consistent with known characteristics of each group. Evaluation is therefore comparative and behavioural, focusing on the separation of group-level patterns and the interpretability of sentence-level evidence. This evaluation approach is consistent with established practices in anomaly detection and exploratory data analysis, where systems are assessed using known exemplars rather than trained to optimise classification accuracy [26, 27]. Similar patterns are observed in intelligence analysis, where accumulated evidence and representative cases guide interpretation in the absence of deterministic decision rules [28].

The outputs of the pipeline should be interpreted as descriptive rather than predictive. Aggregate measures, such as the centroid of sentence-level sentiment scores, provide a summary of the balance of positive and negative signals, while the full distribution captures the presence of mixed or extreme values. As discussed in **sec:aggregation (placeholder)** and **sec:mixed-sentiment (placeholder)**, these representations highlight different aspects of the underlying data and should be considered jointly.

Across the SOI groups, the results exhibit consistent qualitative patterns. Individuals in the “awarded” group tend to show predominantly positive sentiment, while those in the “infamous” group show predominantly negative sentiment. The “forfeited” group occupies an intermediate position, characterised

by mixed sentiment distributions and overlapping centroid values. This overlap reflects the heterogeneous nature of the underlying evidence rather than an absence of signal, and reinforces the limitations of interpreting aggregate measures in isolation.

In this context, the methodology is best understood as a tool for ranking or flagging cases based on the presence of potentially relevant evidence. In particular, the identification of negative sentence-level signals provides a mechanism for surfacing information that may warrant further investigation. As discussed in **sec:minos (placeholder)**, the approach explicitly incorporates asymmetric error costs, prioritising the detection of negative evidence over the avoidance of false positives.

Human evaluation is therefore not an optional safeguard but a fundamental requirement. The system is designed as an evidence-gathering and prioritisation pipeline, in which automated methods structure large volumes of unstructured text and highlight potentially important signals, while the interpretation of those signals—including their relevance, reliability, and context—remains the responsibility of the human analyst.

This design is consistent with established principles for responsible AI systems [29], which emphasise interpretability, reliability, and the retention of human oversight in settings where decisions are context-dependent and cannot be reduced to deterministic rules. In this sense, the methodology does not attempt to replace human judgement, but to augment it by enabling more efficient and systematic analysis of publicly available information at scale.

Overall, the results demonstrate that meaningful sentiment signals can be extracted from noisy open-source data, but that these signals must be interpreted with care. The combination of sentence-level analysis, distributional representation, and conservative detection of negative evidence provides a practical framework for supporting decision-making in settings characterised by limited ground truth and non-deterministic outcomes.

## 5 Conclusions and Future Work

This paper has outlined a data-science approach and system for evaluating the Queen’s Award candidates by analyzing publicly available information about them on the Internet. The system scrapes the Internet for publicly available articles on individuals and then analysing these articles to determine their suitability for the award. The data-pipeline created for the algorithm is highly efficient and about 60 candidates can be evaluated in less than an half an hour. A human sifting through the background of a single individual can potentially take much longer than that. Therefore the system is orders of magnitude more efficient than doing the process manually. However, the developed system assumes that a human analyst or award panel would be responsible for going through key information flagged by the system to come up with the final decision.

A novel sentiment analysis algorithm called MINOS has been proposed by this research which is tailored for the current purpose as compared to previous sentiment analysis algorithms like AFINN and VADER. The detailed results shown in this article justify the need for the various steps in the algorithm and their corresponding benefits. The system created by this work finds application in evaluating the backgrounds of candidates not just for the Queen’s awards but many such honours and awards which expect high standards of conducts from the awardees. A possible future extension to this work would be the periodic evaluation of upcoming candidates as well as previous awardees such that any risks for forfeiture can be flagged early.

## Appendix A Individuals or Subjects of Interest

**Table A1:** The “awarded” group who received an Honours and maintained their award.

SOI	Name	Award	Citation	Gender	Year
A.1	Professor Dame Winifred Mary BEARD OBE	DBE	Study of Classical Civilisation	F	2018
A.2	Professor Sir Sushantha BHATTACHARYYA CBE	KBE	Higher Education and Industry	M	2003
A.3	Sir Donald BRYDON CBE	KBE	Business and charity	M	2018
A.4	Allan COOK	CBE	Defence and Aerospace Industries	M	2018
A.5	Hugh David FACEY MBE	OBE	Manufacturing, Innovation, Exports and Employee Ownership	M	2018
A.6	Rear Admiral Philip Duncan GREENISH	CBE	Military Division	M	2002
A.7	Professor Carole HILLENBRAND OBE	CBE	Understanding of Islamic History	F	2018
A.8	Dr Mohammed Kamal HOSSAIN	OBE	Industry	M	2009
A.9	Professor Sir James HOUGH OBE FRS FRSE	KBE	Detection of Gravitational Waves	M	2018
A.10	Sandra KERR OBE	CBE	Equality and to Diversity	F	2019
A.11	John Nigel Kirkland	OBE	Derbyshire	M	1999
A.12	Professor Richard Ian KITNEY	OBE	IT in Health Care	M	2001
A.13	Ursula Frances Rosamond LIDBETTER	MBE	Business in Lincolnshire	F	2011
A.14	Professor John Neil LOUGHHEAD OBE	CB	Research and Development in the Energy Sector	M	2018
A.15	Miss Maria McCAFFERY MBE	OBE	Renewable Energy Sector	F	2017
A.16	Professor Carol PROPPER CBE FBA	DBE	Economic Policy and Public Health	F	2020
A.17	Dr. Frances Carolyn SAUNDERS CB	DBE	Science and Engineering	F	2018
A.18	Ms Jennifer Margaret SAUNDERS OBE	CBE	Tackling Fuel Poverty	F	2018
A.19	Jack Crossley TORDOFF MBE	OBE	Business and West Yorkshire	M	2018

**Table A2:** The “test” group of individuals who committed serious crime.

SOI	Name	Criminal activity	Gender <sup>1</sup>
T.1	Ajmal Kasav	Terrorism and extremism	M
T.2	Bruce Reynolds	Theft, assault, drug dealing	M
T.3	Charles Sobhraj	Murder, attempted murder	M
T.4	Dale Cregan	Murder	M
T.5	Dennis Nilsen	Murder, attempted murder	M
T.6	Ghislaine Maxwell	Sex trafficking	F
T.7	Harold Shipman	Murder	M

<sup>1</sup>The gender proportion in this table is approximately aligned with UK Criminal Justice System statistics on violent and serious crime, e.g. §8 of [30]

**Table A2:** (continued)

SOI	Name	Criminal activity	Gender
T.8	Harvey Weinstein	Sexual offences	M
T.9	Howard Marks	Drug dealing	M
T.10	Jeffrey Epstein	Sexual offences	M
T.11	John Bodkin Adams	Fraud, perverting the course of justice	M
T.12	Osama Bin Laden	Terrorism and extremism	M
T.13	Oscar Pistorius	Murder	M
T.14	Peter Sutcliffe	Murder, attempted murder	M
T.15	Reginald Kray	Murder, accessory to murder	M
T.16	Robert Maxwell	Not convicted <sup>2</sup>	M
T.17	Ronald Kray	Murder	M
T.18	Samantha Lewthwaite	Not convicted <sup>3</sup>	F
T.19	Samuel Little	Murder, attempted murder	M
T.20	Shamima Begum	Not convicted <sup>4</sup>	F

**Table A3:** The “forfeited” group of recipients who were stripped of the Honour.

SOI	Name	Award	Citation	Gender	Year awarded	Forfeiture reason	Year forfeited
F.1	Anne Ganley	MBE	Employment	F	2012	Perverting the course of justice	2017
F.2	Ashuk Ahmed	MBE	Young people	M	2009	No criminal convictions found	2019
F.3	Craig Martin Burrows	MBE	Charitable and voluntary work	M	2004	Sexual offences	2017
F.4	David John Kemp	MBE	Education	M	2013	Sexual offences	2017
F.5	Derek Charles Eaglestone	MBE	Charitable and voluntary work	M	1994	Sexual offences	2017
F.6	Ian Richard Swingland	OBE	Conservation	M	2006	Fraud	2017
F.7	Ian Strong	MBE	Rural Community in Yorkshire	M	1997	No criminal convictions found	2019
F.8	Jawaid Mohammed Ishaq	MBE	community relations in South Humber-side and North Lincolnshire	M	2000	Fraud	2016
F.9	John Anthony Coatman	MBE	young people	M	2011	Sexual offences	2019
F.10	Malcolm Belchamber	MBE	Littlehampton community	M	2004	Fraud; Forgery and Counterfeiting Act 1981	2017
F.11	Michael Nathan Cohen	MBE	Chorlton Probation Hostel	M	1998	Sexual offences	2018
F.12	Jo Shuter	CBE	Education	F	2010	Professional misconduct	2015
F.13	Patrick Robert John Rock	OBE	Political service	M	1992	Sexual offences	2017

<sup>2</sup>Widespread posthumous evidence of fraud.<sup>3</sup>Warrant issued for charges of possession of explosives and conspiracy to commit a felony.<sup>4</sup>Deprived of UK citizenship due to links to terrorism and extremism.

**Table A3:** (continued)

SOI	Name	Award	Citation	Gender	Year awarded	Forfeiture reason	Year forfeited
F.14	Paul Symonds	OBE	Community Relations in Northern Ireland	M	2007	No criminal convictions found	2017
F.15	Paula Marie Vasco-Knight	CBE	Health services	F	2013	Fraud	2017
F.16	Philip Anthony Knight	OBE	British Honorary Consul-General, Antwerp	M	2001	No criminal convictions found	2017
F.17	Philippa Ann Rodale	MBE	Animal Welfare and to the community in Dorset	F	2007	Professional misconduct; animal welfare charges	2017
F.18	Robert Stanley Poots	MBE	Education	M	2010	Fraud	2017
F.19	Rolf Harris	CBE	Entertainment and the arts	M	2006	Sexual offences	2015
F.20	Trevor George Francis	MBE	Fife community	M	2012	Sexual offences	2017

## References

- [1] Honours and Appointments Secretariat: History. In: The Honours System of the United Kingdom, (2022). <https://honours.cabinetoffice.gov.uk/about/history/>
- [2] Phillips, H.: Review of the Honours System. Crown, ??? (2004)
- [3] Cabinet Office: Having Honours Taken Away (Forfeiture). <https://www.gov.uk/guidance/having-honours-taken-away-forfeiture> (2021)
- [4] Web Graph SIA: Brief History of Web Scraping. <https://webscraper.io/blog/brief-history-of-web-scraping> (2021)
- [5] Hirst, G.J.: Anaphora in natural language understanding: A survey. Master’s thesis, Department of Engineering Physics, Research School of Physical Sciences, The Australian National University (1979)
- [6] Birjali, M., Kasri, M., Beni-Hssane, A.: A comprehensive survey on sentiment analysis: Approaches, challenges and trends. *Knowledge-Based Systems* **226**, 107134 (2021)
- [7] Hutto, C.J., Gilbert, E.: VADER: A parsimonious rule-based model for sentiment analysis of social media text. In: Proceedings of the Eighth International Conference on Weblogs and Social Media, ICWSM 2014 (2014)
- [8] Nielsen, F.A.: A new ANEW: Evaluation of a word list for sentiment analysis in microblogs (2011). Creative Commons Attribution 3.0 Unported
- [9] Storrs, E.: Celebrities on Twitter: Tweet Sentiment Analysis with ULMFiT. Medium, <https://medium.com/@epstorrs/celebrities-on-twitter-tweet-sentiment-analysis-with-ulmfit-3f4746f87b02> (2018)
- [10] Silaparasetty, N.: Twitter Sentiment Analysis for Data Science Using Python in 2022. Medium, <https://medium.com/@nikitasilaparasetty/twitter-sentiment-analysis-for-data-science-using-python-in-2022-6d5e43f6fa6e> (2022)
- [11] Steinberger, R., Hegele, S., Tanev, H., Della Rocca, L.: Large-scale news entity sentiment analysis. In: Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017, pp. 707–715. INCOMA Ltd., Varna, Bulgaria (2017)
- [12] Phillips, H.: Review of the Honours system 2004. Corporate report, Cabinet Office (July 2004)
- [13] Jay, A., Evans, M., Frank, I., Sharpling, D.: I.2 operation of the Honours system. In: Allegations of Child Sexual Abuse Linked to Westminster Investigation Report, (2020)
- [14] Armstrong, H.: Honours: History and reviews. Briefing Paper 02832, House of Commons Library (February 2017)
- [15] Cabinet Office: Operation of the Honours system 2019. Corporate report, Cabinet Office (July 2019)
- [16] Software Freedom Conservancy: The Selenium Browser Automation Project. <https://www.selenium.dev/documentation/> (2004)
- [17] Richardson, L.: Beautiful Soup. <https://www.crummy.com/software/BeautifulSoup/> (2004)
- [18] Margoni, T., Kretschmer, M.: A Deeper Look into the EU Text and Data Mining Exceptions: Harmonisation, Data Ownership, and the Future of Technology. Zenodo (2021). <https://doi.org/10.5281/zenodo.5082012>
- [19] Bergman, J., Popov, O.B.: Exploring dark web crawlers: A systematic literature review of dark web crawlers and their implementation. *IEEE Access* **11**, 35914–35933 (2023)
- [20] Loper, E., Bird, S.: NLTK: The natural language toolkit. Proceedings of the ACL Workshop on Effective Tools and Methodologies for Teaching Natural Language Processing and Computational

- [21] Kiss, T., Strunk, J.: Unsupervised multilingual sentence boundary detection. *Computational Linguistics* **32**(4), 485–525 (2006) <https://doi.org/10.1162/coli.2006.32.4.485>
- [22] Qi, P., Dozat, T., Zhang, Y., Manning, C.D.: Universal dependency parsing from scratch. In: *Proceedings of the CoNLL 2018 Shared Task: Multilingual Parsing from Raw Text to Universal Dependencies*, pp. 160–170. Association for Computational Linguistics, Brussels, Belgium (2018)
- [23] Nielsen, F.A.: AFINN: A New Word List for Sentiment Analysis on Twitter. <https://finnaarupnielsen.wordpress.com/2011/03/16/afinn-a-new-word-list-for-sentiment-analysis/> (2011)
- [24] Nielsen, F.A.: `fnielsen/afinn`. GitHub repository, <https://github.com/fnielsen/afinn/tree/master> (2022)
- [25] Hutto, C.J.: `cjhutto/vaderSentiment`. GitHub repository, <https://github.com/cjhutto/vaderSentiment> (2022)
- [26] Chandola, V., Banerjee, A., Kumar, V.: Anomaly detection: A survey. *ACM Computing Surveys* **41**(3), 15 (2009)
- [27] Tukey, J.W.: *Exploratory Data Analysis*. Addison-Wesley, ??? (1977)
- [28] Heuer, R.J.: *Psychology of Intelligence Analysis*. CIA Center for the Study of Intelligence, ??? (1999)
- [29] Government Digital Service: *Artificial Intelligence Playbook for the UK Government*. <https://www.gov.uk/government/publications/ai-playbook-for-the-uk-government>. Accessed 2026 (2025)
- [30] Ministry of Justice: *Women and the Criminal Justice System 2021*. <https://www.gov.uk/government/statistics/women-and-the-criminal-justice-system-2021/women-and-the-criminal-justice-system-2021#offence-analysis> (2022)