# Detecting patterns in North Korean military provocations: what machine-learning tells us

Taehee Whang[1], Michael Lammbrau[2] and Hyung-min Joo[3]

[1]*Department of Political Science & International Studies, Seoul, South Korea;* [2]*Department of Intelligence Studies, Mercyhurst University, Erie, USA;* [3]*Department of Political Science & International Relations, Seoul, South Korea*
*Email: hjoo@korea.ac.kr*

## Abstract

For the past two decades, North Korea has made a series of military provocations, destabilizing the regional security of East Asia. In particular, Pyongyang has launched several conventional attacks on South Korea. Although these attacks seem unpredictable and random, we attempt in this article to find some patterns in North Korean provocations. To this end, we employ a machine-learning technique to analyze news articles of the Korean Central News Agency (KCNA) from 1997 to 2013. Based on five key words ('years,' 'signed,' 'assembly,' 'June,' and 'Japanese'), our model identifies North Korean provocations with 82% accuracy. Further investigation into these attack words and the contexts in which they appear produces significant insights into the ways in which we can detect North Korean provocations.

# 1 Introduction

Since the conclusion of the Korean War (1950–53), there has been significant violence along the Demilitarized Zone (DMZ). In addition, there have been multiple military clashes at sea and elsewhere along the coast of South Korea. What is the main motivation of Pyongyang to initiate these attacks? Is it possible to analyze North Korean military provocations systematically? In answering these questions, a primary obstacle has been the relative dearth of information regarding the intentions of Pyongyang. As one put it, North Korea has been 'the longest running intelligence failure' (Litwak, 2007) or 'North Korea could have been on Mars for [the outside world] knew about it. It was a faraway land of unknowns and unknowables explored mostly by space probes, and in this case, spy satellites' (Sigal, 1998).

Existing theories of IR fail to provide an adequate guide to understand North Korea military provocations. For instance, consider realism. According to Mearsheimer (2001), three factors affect power calculations of a country: possession of nuclear forces, separation by large bodies of water, and a power distribution. Among them, the first two cannot explain variations in North Korean military provocations because Pyongyang has the upper hand in nuclear capability against Seoul and a territorial proximity between the two Koreas is a constant. By contrast, the distribution of power can influence the extent of fear that North Korea may have because of increasing power asymmetry since the collapse of the Soviet bloc. What is unclear, however, is the level of resolve North Korean has to initiate military attacks. Without data to estimate how willing Pyongyang is to use force, realism is limited in explaining how profound North Korean fears are and, hence, why Pyongyang resorts to violent military attacks occasionally.

To cope with the paucity of reliable information, two lines of research have emerged in North Korean studies: a survey/interview of North Korean refugees and an analysis of North Korean newspapers. About the latter, the Korea Central News Agency (KCNA) has published government-approved articles since 1996. In recent years, several scholars have focused on the KCNA in hopes of distilling useful insight from it (McEachern, 2010; Rich 2012a,b; Joo, 2014). Based on counting the frequency with which particular terms or phrases appear in KCNA news articles, these works have yielded some interesting findings

about the linguistic features of KCNA articles and their correlation with nuclear policies, political rhetoric, economic trends, and social changes of North Korea.

By contrast, this article embarks on a new approach: a text-classification approach based on a supervised machine-learning technique. In particular, we develop a model that can distinguish the period of imminent North Korean provocations from peace time, by using the KCNA as our data and supervised machine-learning as our method. For our cases, we select all five North Korean attacks between 1997 and 2013, (i) the First Battle of Yŏnpyŏng on 15 June 1999, (ii) the Second Battle of Yŏnpyŏng on 29 June 2002, (iii) the Battle of Daechŏng on 10 November 2009, (iv) the sinking of the Cheonan naval ship on 26 March 2010, (v) the shelling at Yŏnpyŏng Island on 23 November 2010. Each of these incidents is a conventional North Korean attack resulting in more than one casualty on one or both sides.

Our analysis shows that immediately prior to North Korean attacks, Pyongyang tends to increase its use of five key terms in the KCNA: 'years', 'signed', 'assembly', 'June', and 'Japanese.' Our investigation into the contexts in which they appeared in KCNA articles shows that Pyongyang often employs terms like 'June', 'years', and 'Japanese' to nostalgically invoke past battles of Kim Il-sung against Japanese colonialism. Moreover, the term 'signed' indicates that the KCNA quotes 'official commentaries' published in *Rodong Sinmun*, the mouthpiece of the ruling Workers' Party of Korea (WPK), right before an attack. Finally, a social network analysis of key terms shows that the word 'assembly' refers to the Supreme People's Assembly (SPA). The high correlation of the term 'assembly' with military attacks allows us to conjecture that provocations are often premeditated insofar as they are preceded by increased SPA activities. These findings provide us with a basis for further research into North Korean military provocations. By differentiating threat articles from non-threat items, our model can serve as a useful indicator for imminent North Korean aggressions.

## 1.1 Logistics: data and cases

### 1.1.1 Data: KCNA

As the sole news agency of North Korea, the KCNA provides daily reports of North Korean newspapers (e.g. *Rodong Sinmun* of the ruling WPK and *Minju Chosun* of the North Korean government), television broadcasts (e.g. Korean Central Television), and radio broadcasts (e.g. the Korean Central Broadcasting System). Since 1996, the KCNA has published daily news articles via its server in Japan (www.kcna.co.jp). On a typical day, the KCNA publishes 20–40 articles, including reports on activities of the ruling North Korean elite, official statements of the North Korean government (e.g. an official statement from the Ministry of Foreign Affairs on nuclear issues), several articles selected from *Rodong Sinmun* and *Minju Chosun*, miscellaneous news about North Korean society, and reports of recent developments in foreign countries.

Scholars working on North Korea have relied on its two medias: *Rodong Sinmun* and the KCNA. As the newspaper of the WPK, *Rodong Sinmun* is regarded as the official mouthpiece of the North Korean regime. As a result, it has become popular for scholars, especially in South Korea, to conduct content analyses of *Rodong Sinmun* in order to identify trends or policy shifts of the North Korean government (Koh, 2013). From our viewpoint, however, *Rodong Sinmun* has two weaknesses. First, it is inappropriate for our project because its primary target is the domestic audience of North Korea (thus, published only in Korean). Given that *Rodong Sinmun* is written for a domestic audience, it is not a proper place to look for signs, patterns, or indicators of Pyongyang that its relations with the outside world (especially, the U.S. and South Korea) are at a breaking point and that a military conflict of some sort is about to occur. Instead, the KCNA with its focus on foreign audiences (thus, published in four different languages – English, Spanish, Russian, and Korean) provides a better source to detect such signs or patterns. Second, the KCNA provides a better dataset from a technical viewpoint. Although North Korea has made *Rodong Sinmun* available on internet (www.kcna.co.jp/today-rodong/rodong.htm) in recent years, only few selected articles after 2002 are

available while only titles are provided for the rest of *Rodong Sinmun* articles. By contrast, all the articles in the KCNA after 1996 are available on the internet, thus providing better data for our machine-based text-classification analysis to detect any signs, patterns, or indicators from Pyongyang that a military strike is likely to occur. As a result, the KCNA is used as the main source of data.

Although our case selection of conventional North Korean provocations begins in 1997 because KCNA data is available after that year, it is more than a technical convenience to use the year 1997 as a starting point. It also overlaps with the succession process from Kim Il-sung to Kim Jong-il. When Kim Il-sung passed away on 8 July 1994, Kim Jong-il took the traditional three-year mourning period (1994–96) before he assumed the official title of General Secretary of WPK in 1997 to rule the country. As a result, the year 1997 serves as a starting point in our project not only for a technical reason (i.e. the availability of the KCNA dataset after 1996) but also for a substantive reason (i.e. it included North Korean military provocations in the post-Kim Il-sung era). In particular, Pyongyang has launched five conventional military strikes in the post-Kim Il-sung period.

## 1.2 Cases: five conventional military crises since 1997

Figure 1 shows five North Korean conventional military attacks between 1996 and 2013: (i) the First Battle of Yŏnpyŏng on 15 June 1999; (ii) the Second Battle of Yŏnpyŏng on 29 June 2002; (iii) the Battle of Daechŏng on 10 November 2009; (iv) the sinking of the Cheonan naval ship on 26 March 2010; and(v) the shelling at Yŏnpyŏng Island on 23 November 2010. For our case selection, three criteria are used. First, each of these attacks caused one or more casualties on at least one side. Second, all five attacks used conventional non-nuclear weapons.[1] Third, all five attacks were initiated by North

---

1   In this article, we have excluded all events associated with the North Korean nuclear crisis. In a separate project, however, we employ a machine-learning technique to develop to detect significant signs or patterns of Pyongyang that it is about to conduct a nuclear test. Preliminary research shows an interesting contrast between conventional provocations and nuclear crisis in the North Korean case. As will be shown, a single platform (covering the entire period from Kim Jong-Il to the Kim Jong-Un) outperforms a double platform (one model for the KJI period and another model for the KJU era) in the case of North Korean conventional provocations. Simply put, there has not been a major policy shift in
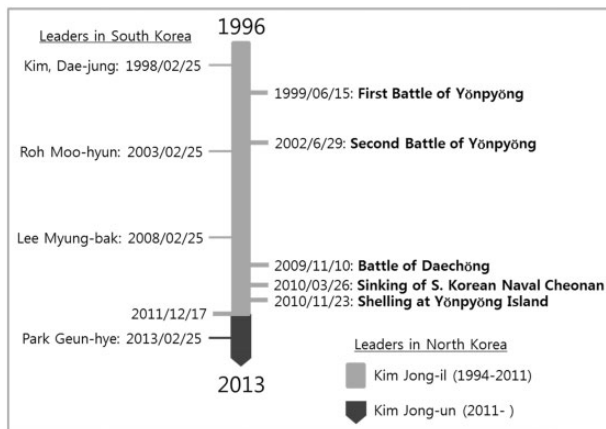
**Figure 1** Five conventional North Korean military attacks

Korea. Below is a brief description of the five North Korean military provocations.

### 1.2.1 The first battle of Yŏnpyŏng:15June1999

In 1999, Pyongyang claimed that South Korea trespassed the Northern Limit Line (NLL) in the Yellow Sea. On 7 June 1999, when North Korean patrol ships and fishing boats crossed the NLL, the South Korean navy responded by increasing its patrol. After a few days of continual trespassing, the South Korean government issued a warning and deployed two patrol corvettes (Hong, 2012). When Pyongyang continued to ignore warnings, the South Korean navy blocked North Korean boats, ramming them (Macfie, 2013). After a few skirmishes, the main battle occurred on 15 June 1999 when four North Korean patrol ships trespassed across the NLL, soon joined by three North Korean torpedo boats and three additional patrol ships. With a total of ten battleships, North Koreans launched a 25 mm cannon shell

Pyongyang as far as its conventional provocations are concerned. By contrast, a double platform outperforms a single platform in the case of the North Korean nuclear crisis, indicating a major policy shift from Kim Jong-Il to Kim Jong-Un. At the moment, we are trying to find out whether such a shift indicates an increasing intention of the new North Korean leadership under Kim Jong-Un to proclaim a 'nuclear power' status by developing its nuclear program further, instead of negotiating over it as his father had done on several occasions.

(Park, 2009). In response, the South Korean navy fired with 40 mm and 76 mm machine guns. When the battle was over, a North Korean torpedo boat sank, a large patrol ship crashed, and four patrol ships sustained damage. In the process, approximately 30 North Korean soldiers were killed and 70 were wounded. As for South Korea, four patrol killers and one patrol corvette were damaged with nine soldiers wounded (Moore and Hutchison, 2010).

### 1.2.2 The second battle of Yŏnpyŏng: 29 June 2002

By 2002, North Korean ships frequently crossed the NLL, only to be chased back by South Korean patrol vessels. On 29 June 2002, two North Korean patrol ships crossed the border, ignoring warnings from South Korean navy speedboats. At 10:25 am, the two North Korean patrol boats attacked nearby South Korean vessels, with its 85 mm gun, 25 mm auxiliary gun, and hand carried rockets. In response, the South Korean patrol ships returned fire (Sohn, 2002). A 20-minute battle resulted in the death of four South Korean marines, one missing, and 18 wounded. On the North Korean side, approximately 30 sailors were killed or injured. While South Korean vessels were partly damaged, one of the North Korean vessels was towed away in flames (Global Security, 2002).

### 1.2.3 The battle of daechŏng: 10 November 2009

On 10 November 2009, a North Korean patrol vessel trespassed the NLL. Soon after, two 130-ton South Korean vessels issued warnings but the North Korean vessel ignored them. When the South Korean ships fired warning shots, the North Korean vessel began firing, leading to a 2-minute battle near Daechŏng Island, located 18 miles off the North Korean coast. The North Korean patrol vessel also attacked a South Korean high-speed patrol vessel (Kim, 2009). In response, the South Korean vessel countered with approximately 200 shots. When the battle was over, there were no South Korean casualties, but North Korea suffered one casualty and three injuries with its naval vessel 'wrapped in flames' (Choe, 2009).

### *1.2.4 Sinking of the cheonan naval ship: 26 March 2010*

On 26 March 2010, the South Korean naval ship Cheonan sank into the Yellow Sea. At 9:22 pm, an explosion occurred in the 1,200-ton warship that was sailing by Baengnyŏng Island where the two Koreas had clashed numerous times before (Cha, 2010). In the end, 46 lives were lost and it was quickly suspected that the ship 'had been hit by an external 'non-contact' explosion', with North Korea as the prime suspect (Sudworth, 2010). Pyongyang denied its involvement and claimed that the incident had been contrived by South Korea. A six-week-long investigation by international experts proved the involvement of North Korea in the attack.

### *1.2.5 Shelling of Yŏnpyŏng island: 23 November 2010*

On 23 November 2010, North Korea fired artillery shells at Yŏnpyŏng Island near the NLL. About 200 shells hit the island and set fire to dozens of buildings. The barrage killed two South Korean citizens and two marines, while injuring three civilians and 17 soldiers. The attack began when South Korea was practicing military drills near the NLL despite North Korean warnings. North Korea fired three separate barrages, with dozens of artillery shells in each barrage. In return, South Korea responded by firing 80 rounds from K9 155 mm self-propelled Howitzers (Kim and Kim, 2011). When the battle was over, more than 50 civilian homes were in flames.

## 1.3 Research method: supervised machine learning

Although the North Korean nuclear crisis has been the focus of the international community in recent years, the history of North Korean military provocations dates much further back. For instance, Pyongyang destabilized an already precarious security environment in the Korean Peninsula, by launching a series of provocations such as the hijacking of a South Korean airline (1958), the Korean DMZ Conflict known as 'the Second Korean War' with more than 700 casualties (1969), the hijacking of the USS Pueblo (1968), the notorious Axe Murder Incident (1976), the bombing of the Korean Air Line 858 (1987), and so on.

Not surprisingly, many scholars have attempted to analyze North Korean military provocations, trying to identify their 'causes' (Jung, 2013; Lee, 2014; Ko, 2015), main 'goals' in those attacks (Jung, 2008; Kang, 2013), proper 'policies' to curtail further threats (Oh, 2011; Kim, 2012), and so on. Despite sincere efforts, earlier studies suffered from one notorious problem: a lack of reliable data. As one put it, North Korea has been 'the blackest of black holes' (Litwak, 2007, 289). Under such circumstances, scholars had little choice but to make educated guesses – 'guesstimations' – regarding unknown intentions, goals, or likely moves of the North Korean regime. As a result, the existing literature on North Korean military provocations has been driven less by reliable data but more by subjective interpretations.

By contrast, we rely on the official North Korean media KCNA. As for the method, we use a supervised machine-learning technology to maximize the use of the KCNA dataset that covers various aspects of North Korea from 1997 to 2013. The main advantage of our method comes from the quality of the data; that is, our findings are data-driven and thus more objective than previous works relying on subjective interpretation of selected observations. Given the paucity of reliable information on North Korea, it is important to analyze the content of official KCNA articles that are publicly available. A supervised machine-learning technology is the optimal method to process such data objectively.

Supervised machine-learning consists of three steps:(i) data collection and document labelling, (ii) pre-processing, and (iii) model extraction and analysis.[2] First, we obtained our dataset from the KCNA website (http://www.kcna.co.jp). Since there are five North Korean military attacks in our case, we pulled 10 tranches of articles from all the KCNA articles published from 1997 to 2013. We then labelled five of these tranches 'threat' while labelling the other five as 'non-threat.' All

---

2   We tried a variety of machine learning algorithms, such as Random Forests, Support Vector Machines, and Conditional Inference Trees. Also, we cross-validated their results. Our findings demonstrate that these algorithms produced similar results. Moreover, all of these algorithms identified the similar pattern-detecting terms for classifying the KCNA articles. In this article, we reported the results based on the Conditional Inference Tree algorithm, which ran the tree-structured regression models through binary recursive partitioning in a conditional inference framework. For more details on the Conditional Inference Tree algorithm, please see the R Package 'Party' for Conditional Inference Trees ('ctree') at http://cran.r-project.org/web/packages/party/party.pdf.
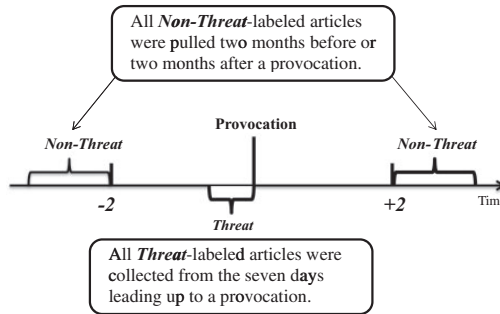
**Figure 2** Labelling threat and non-threat articles

the articles published within a week preceding each attack is labelled as 'threat' and the rest of KCNA articles were treated as 'non-threat.' In total, there were 1,624 KCNA articles published during this period, with 487 labelled as 'threat' and 1,137 labelled as 'non-threat.' For the threat articles, we extracted all the articles published within a week of each North Korean military provocation (Fig. 2). For the non-threat articles, we selected tranches of news articles for a randomly chosen 10-day period from at least two months before or after a North Korean attack. In so doing, our goal was to capture articles that were unrelated to a North Korean attack, thus labelled as 'non-threat.'

Second, the selected dataset is then 'preprocessed.' Since the original articles are not ready for an automated text analysis, a cleaning process called 'preprocessing' is necessary to prepare them for machine-learning. Pre-processing is the standard procedure before the application of an automated text analysis. Following the standard preprocessing procedure, we remove all numbers and stopwords (e.g. a, the, in, to, etc.).[3] We then transformed the body of the text into a document term matrix that was composed of rows of news articles followed by columns of terms. The document term matrix turned preprocessed texts

---

3  In this article, we use one of the most popular preprocessing methods known as the 'bag of words' concept (Jurafsky and Martin, 2009). It discards the use of word order as a factor and removes the so-called 'stopwords' from the data. The stopwords include punctuation, capitalization, common words (e.g., at, an, the, etc.), and numbers. For a full list of the stopwords used in our research, please visit http://jmlr.csail.mit.edu/papers/volume5/lewis04a/a11-smart-stop-list/english.stop. Although preprocessing reduces the amount of information, it has been shown by researchers that a simplification of text via preprocessing is sufficient to infer valuable models (Hopkins and King, 2010).
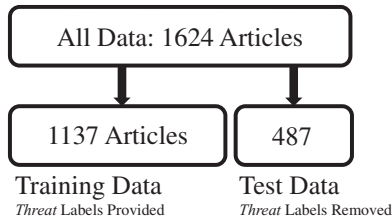
**Figure 3** Training and test data sets

into quantifiable values (e.g. term counts and frequencies) for each cor-
responding article.

Finally, once we had collected and preprocessed the KCNA data, we
split the whole data into two subsets: a training dataset and a test
dataset. From the complete set of articles from both threat and non-
threat tranches (1,624 articles in total), we randomly selected 70%
(1,137 articles) for the training dataset (Fig. 3). The labels 'threat' or
'non-threat' for each article were included in the training dataset for
the purpose of automated machine-learning: that is, to develop a
model that can select pattern-detecting features based on a priori clas-
sifications. Basically, the key pattern was frequency of occurrence. The
frequency with which certain words and short phrases appeared in
threat or non-threat articles determined how they were selected and
weighted as pattern-detecting features in the model. The trained
model was then applied to the remaining test dataset, which was used
solely for testing the accuracy of our model. Importantly, we removed
all the threat and non-threat labels from each article in the test
dataset. The purpose for doing so was to challenge our model to use
what it had learned (from the training dataset) about significant fea-
tures (i.e. a term frequency rate) to accurately classify KCNA articles
(in the test dataset) as either a threat or a non-threat, without relying
on labels.[4]

---

4    For replication, our data is available at https://dataverse.harvard.edu/dataset.xhtml?
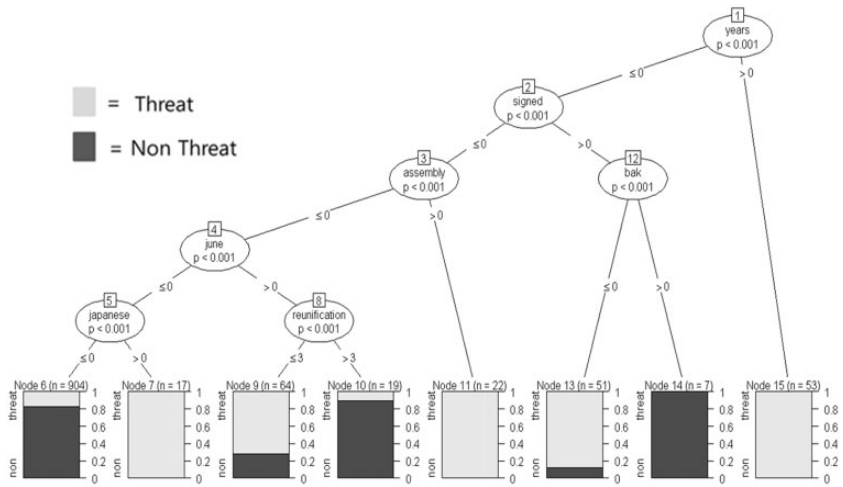     persistentId=doi%3A10.7910%2FDVN%2FB8CWWD.

**Figure 4** First conditional inference tree

# 2 Results

## 2.1 Training data results

Figure 4 shows the model we obtained from the training dataset using a supervised machine-learning. The key terms distinguishing North Korean threats from a peaceful period were 'years', 'signed', 'assembly', 'June', and 'Japanese.' When certain conditions were satisfied (as explained below), the appearance of these terms in KCNA articles could play the role of a pattern-indicator that a North Korean military threat was on the horizon.

According to our model, the first and strongest indicator of a North Korean military provocation is the appearance of the term 'years.' In the training dataset, all KCNA articles in which the word 'years' appeared at least once (53 articles in total) were published within a week of a North Korean military strike. From our viewpoint, the most reliable sign of an imminent North Korean military provocation is a sudden increase in the frequency of the term 'years' in KCNA articles.

When the term 'years' is absent, the second best indicator of a North Korean military attack is the word 'signed.' Approximately 85% of KCNA articles (51 articles in total) that had the term 'signed' without the word 'bak' (from former South Korean president Lee Myung-bak)

were published within a week of a North Korea military provocation. As a result, a sudden spike in the use of the word 'signed' (without 'bak') indicates that a North Korean military provocation is around the corner. In this respect, however, there is a further condition that must be satisfied. When 'signed' and 'bak' appear together in the same article, they indicate a pattern of a peaceful situation instead. As a result, the pattern-detecting role of the term 'signed' appears to be contingent upon the presence or absence of another term 'bak.'

According to our model, the third index indicating that Pyongyang may be preparing for a military operation is the term 'assembly.' In fact, all the KCNA articles that included the word 'assembly' without 'years' or 'signed' (22 articles in the training dataset) were published within one week of a North Korean military attack. As a result, a sudden increase in the frequency of the word 'assembly' may be a good pattern indicator that a North Korean military strike is on the horizon, even in the absence of stronger threat terms such as 'years' and 'signed'.

The fourth indicator of a North Korean military threat is the term 'June'. Unlike other indicators, however, the role of 'June' is conditional. In the absence of other stronger signs (i.e. years, signed, and assembly), the term 'June' can play the role of a significant indicator of a North Korean military attack only if another term 'reunification' appears once or less in the same article. To our surprise, however, the term 'June' turns into a strong indicator of a *non*-threat situation if the word 'reunification' appears more than twice in the same article. As a result, it seems that the word 'June' implies an increasing North Korean threat only if it is not closely associated with another key pattern indicator, 'reunification'.

The last index of a North Korean military provocation is the word 'Japanese'. When other (and stronger) terms such as 'years', 'signed', 'assembly', and 'June' were not present, all the KCNA articles that included the term 'Japanese' (17 articles in the training dataset) appeared within one week of a North Korean military provocation. As a result, the use of the term 'Japanese' serves as a good pattern indication that, in the absence of other attack words, Pyongyang is moving dangerously close to a military strike.

**Table 1** Model accuracy against test dataset

| Model | Actual Threat | Non-threat | | |
|---|---|---|---|---|
| Threat | 74 | 13 | Positive predictive value | 0.85 |
| Non-threat | 73 | 327 | Negative predictive value | 0.82 |
| | Sensitivity | Specificity | Overall accuracy | |
| | 0.50 | 0.96 | 0.82 | |

## 2.2 Testing the model

There are 904 KCNA articles that do not include any of the indicators discussed above. Because they do not contain any pattern indicator of a North Korean attack, our model identifies them as non-threats. In reality, however, about 20% of these articles were published within one week of the five military provocations by Pyongyang. As a result, our model has roughly 80% accuracy in identifying North Korean military threats. What is encouraging is that our model has identified five key pattern indicators of increasing North Korean threats. When these key terms are put together under certain conditions, they correctly identify 80% of articles in the training dataset as threat articles or non-threat items. Put differently, the model can accurately classify in 8 of 10 cases whether various messages from Pyongyang are real threats or just rhetoric.

Although 80% accuracy is impressive, it is too early to be optimistic. After all, 80% accuracy was achieved with the training dataset, from which our model was originally derived. If we were impressed by its 80% accuracy rate, we would resemble a case study specialist who developed an elaborate theory from a few cases and then applied it back to those original cases only to be impressed by how accurate his/her theory was. The real test consists of applying our model to cases *other than* those from which it is derived. This is the reason why we divided the entire KCNA data into two subsets: the training dataset from which our model is developed and the test dataset to which it will be applied as a *real* test.

Table 1 shows the result of our model against the test dataset. Numbers in the cells indicate whether there was agreement or discrepancy between model predictions and actual cases. For example, our

model classified 327 cases (the lower right cell) as non-threat and those articles were, in fact, published during non-threat weeks. In addition, the model classified 74 articles (the upper left cell) as threats and those articles were actually published during threat weeks. When we apply our model to the test dataset, its overall accuracy turns out 82.3%, roughly the same level we obtained with the training dataset. Although our model is deduced from the training dataset, it does not lose categorizing capacity when applied to the test dataset. Specifically, of 487 articles in the test dataset, our model correctly classified 401 articles (82.3%) as either threats or non-threats. The model is very effective in identifying harmless rhetoric from Pyongyang: that is, messages not resulting in actual military attacks. When applied to a total of 340 non-threat articles in the test dataset, our model successfully classified 327 items as noise, with only 13 misses (i.e. 96.7% accuracy). As a result, it is very effective at filtering noise from Pyongyang. By contrast, the pattern-detecting accuracy of our model drops somewhat with respect to threats. Of 147 threat articles in the test dataset, it correctly identifies 74 items as threats while misrecognizing 73 threats as peaceful articles (i.e. 50% accuracy).

## 3 Discussion

What does the model tell us in plain terms? In particular, what are the meanings of the key indicators for a North Korean military threat (i.e. years, signed, assembly, June, and Japanese)? To answer, it is necessary to go back to the original KCNA documents in order to understand the contexts in which these terms were used. A close reading of the KCNA articles that included the five key indicators or 'attack words' reveals several interesting patterns.

First, the North Korean regime tends to emphasize its *history* of military struggle against foreign enemies before it launches armed provocations. It is then understandable that key terms such as 'years' and 'Japanese' often appear in KCNA articles immediately before a military strike. Prior to the second naval clash of Yŏnpyŏng on 29 June 2002, for instance, the North Korean government published a series of articles in which it boasted of the 'years' of North Korea's victorious struggles against foreign enemies, including 'Japanese' colonialism (1910–45). It is possible that Pyongyang invoked the military legacy of

its supreme leader Kim Il-sung in anticipation of an impending conflict, either to boost the morale of its people or to show the outside world that it was determined not to back down.

Second, the most perplexing term in Figure 4 is 'June' because it can be indicative of both threat and peace. Specifically, the word 'June' is an indicator of peace when associated with 'reunification,' but it can also be a sign of a military threat when it does not appear in conjunction with 'reunification.' A careful reading of the KCNA articles that include the term 'June' suggests an explanation for this strange phenomenon. It turns out that the word 'June' can refer to either the Korean War (which broke out on 25 *June* 1950) or the historic summit between Kim Dae-jung and Kim Il-sung on 15 *June* 2000 (which was praised as a major cornerstone for future 'reunification' in the official North Korean press). As a result, when the word 'June' appears in close association with 'reunification,' it refers to the 2000 summit (or 'the June 15 Summit' as it is called in North Korea), thus indicating the peaceful intentions of Pyongyang. By contrast, when the word 'June' appears alone without connection to 'reunification,' it refers to the Korean War, thus conveying a more hostile mood.

Third, it is also interesting to note that the North Korean government tends to publish many 'signed' commentaries from *Rodong Sinmun* in the KCNA immediately before it launches military provocations. In these 'signed' commentaries of *Rodong Sinmun*, Pyongyang usually criticizes foreign countries (especially the United States, Japan, and South Korea) in very harsh terms on various issues. Accordingly, the third term 'signed' seems to correspond to the fact that *Rodong Sinmun*, the official mouthpiece of the North Korean regime, barks very loudly before it actually bites its intended target. A sudden increase in the number of 'signed' commentaries of *Rodong Sinmun* in KCNA articles appears to be a reliable sign that the North Korean government may be turning to an attack mode.

Finally, the last indicator of an imminent threat, 'assembly,' is the easiest term to identify but the hardest one to interpret. As shown in Figure 5, the term 'assembly' is primarily used in reference to the SPA. Because the SPA is the highest North Korean authority with regard to its relations with foreign countries, it is tempting to interpret a sudden increase in references to the SPA prior to military provocations as indicating that Pyongyang is attempting to openly signal its hostile
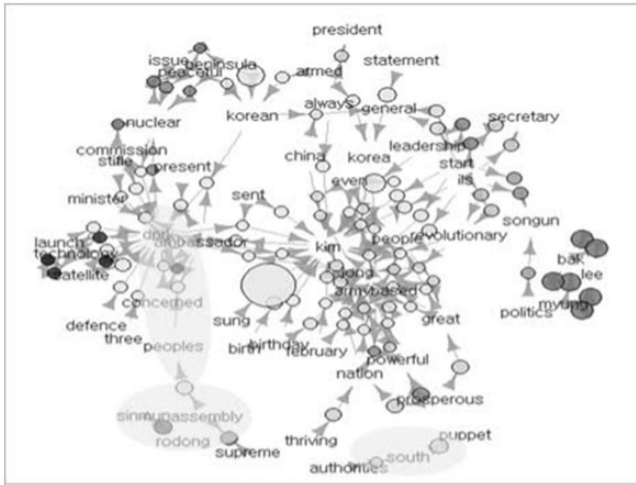
**Figure 5** Social network analysis for key words in KCNA articles

intentions to the outside world. In other words, belligerent messages emanating from the SPA (and reported in KCNA articles) may reveal the determination of the North Korean regime.

Although plausible, the problem with such an interpretation is that a close reading of several KCNA articles using the term 'assembly' shows that messages from the SPA are often *not* dire at all. For example, consider the following article published on 17 November 2010, just days before the North Korean shelling at Yŏnpyŏng Island. 'Kim Yong Nam, President of the Presidium of the DPRK SPA, sent a message of greetings to Qaboos Bin Said, Sultan of Oman, on Wednesday on the occasion of its national holiday' (KCNA 17 November 2010). As this example illustrates, a typical KCNA article with the word 'assembly' describes rather routine business of the SPA, such as how it sent a message of congratulations to a foreign country, how it greeted visiting dignitaries from foreign countries, and so on. Clearly, such mundane messages cannot be a meaningful sign of an imminent North Korean military provocation. While the publication of these articles (containing the term 'assembly') might possibly correspond to an uptick in Pyongyang's diplomatic efforts to strengthen its ties with foreign countries and increase international support prior to a military attack, we need further analyses to substantiate such a hypothesis. At the same time, it also seems clear from the test that the term 'assembly' does not

appear randomly and is somehow linked to the timing of North Korean military provocations. Further research is necessary to make a proper interpretation of this seemingly irrelevant, yet apparently significant, indicator of North Korean military threats.

### 3.1 Robustness checks

Before we conclude, it is necessary to address five issues regarding the robustness of our findings: (i) the rare-event issue (that is, due to the rare nature of a North Korean military attack, our sampling ratio of 7:10 should be justified); (ii) an out-of-sample alternative (that is, instead of building our model by using 70% of data from the five North Korean military strikes and then applying it to the remaining 30% of the data, it may be better to develop a model from the first three strikes and then apply it to the remaining two North Korean provocations); and (iii) a North Korean leadership change effect (that is, instead of elaborating a single model covering 1997 to 2013, it may make sense to develop two separate models [one for the Kim Jong-il period and the other for the Kim Jong-un era] in order to check whether there is a leadership change effect); (iv) a South Korean leadership change effect (that is, it may be better to elaborate two different models [one for the 'conservative' period during Lee Myung-bak and Park Geun-hye, and the other for the 'progressive' era during Kim Dae-jung and Roh Moo-hyun] in order to investigate whether North Korea responded differently to leadership changes in South Korea; and (v) a no-Cheonan alternative (that is, it is worth elaborating a model while excluding the Cheonan naval ship case, because, unlike the remaining four attacks, Pyongyang has denied its involvement in sinking the Cheonan. In this section, these five issues are addressed to check the robustness of our model.

First, there is a discrepancy in the ratio of threat days vs. non-threat days between the data used in our analysis and the actual frequency. As explained earlier, we have defined the threat articles as those published one week prior to each actual crisis, whereas the non-threat items are defined as those chosen randomly for a 10-day period at least two months before or after actual provocations. For all five crises, the ratio of our data is then 35:50 (in days) or 7:10. In reality, however, the number of peace days (i.e. 18 years minus 35 days) is much larger than

the number of crisis days (35 days). When we count all of the peace days that are not included in our data (i.e. when we take all True Negatives into our data), the North Korean provocations become extremely rare events. As a result, a possible criticism is that our approach does not represent the true data generating process in a sampling procedure. While military provocations occur rarely in reality, our sampling procedure exaggerates their likelihood by significantly reducing the number of peace days in the data. Put more technically, our model reduces the number of False Positives while increasing the number of False Negatives.

Despite the criticism, our sampling strategy can be justified for two reasons. First, the rare-event problem is not new in political science. For example, King and Zeng (2001a,b) addressed the same problem in statistical analysis (e.g. logit analysis with rare events). A commonly suggested solution is the 'choice-based or endogenous stratified sampling' in econometrics and the 'case-control design' in epidemiology (Breslow, 1996). The idea is to 'select within categories of Y' such that all crisis days are sampled while we use a small random fraction of non-crisis days. In this respect, our sampling strategy is consistent with such an approach in that we have chosen all crisis days (35) and a random fraction of peace days (50). Not surprisingly, such an approach is commonly used in supervised machine-learning as well. The rare-event issue, also known as a class imbalance problem, is an ongoing issue in supervised machine-learning because the size of one class is often much larger than that of the other class (e.g. premature births, violent civil conflicts, fraudulent credit card transactions, etc.). In supervised machine-learning, two solutions have been suggested: a data sampling technique and an algorithmic modification technique. Whereas data sampling is to under-sample the majority class while over-sampling the minority class, a modification technique is to combine both over- and under-sampling methods at the same time. In this article, we have adopted a data sampling technique in order to address the class imbalance between threat days (a minority class) and non-threat days (a majority class) in North Korean military provocations.

Second, although reducing the number of non-crisis days in our data (under-sampling) while using all the crisis days (over-sampling) is consistent with both statistical and machine-learning literature, there still remains a question. How far should we reduce the number of

peace days? Will a ratio of 1:2, 1:3, or an even more skewed ratio generate a better performance than our model using a 7:10 ratio? To answer this question, we pushed the ratio until we saw a clear sign of the model becoming too skewed towards the majority class (peace days). In our case, it occurred when the ratio between the two classes exceeded 1:3. On this subject, literature on machine-learning clearly shows that the ratio of crisis and non-crisis should not be significantly unbalanced. According to Ertekin (2013), '[i]n problems where the prevalence of classes is imbalanced, it is necessary to prevent the resultant model from being skewed towards the majority (negative) class and to ensure that the model is capable of reflecting the true nature of the minority (positive) class.' Otherwise, the generated model can suffer from an over-fitting problem. For instance, in the face of an extremely skewed dataset (e.g. the *real* ratio of our data would be 35 crisis days vs. 6,535 peace days), an automated machine-learning process naturally becomes 'greedy'; that is, in order to increase its overall accuracy, it increasingly focuses on the larger category (6,535 days) while virtually ignoring the smaller category (35 days). Generating the model in this way would be problematic, however, because our object is to focus on the *minority* category of crisis days, which is less than 0.5% of all days from 1997 to 2013. After all, we are trying to classify upcoming North Korean *attacks*, not peaceful days.

With this goal in mind, we ran robustness checks to see if alternative models yield a better explanatory power when we increased the number of peace days vis-à-vis threat days. Because the maximum limit of an unbalanced ratio for automated machine-learning is 1:3, we attempted both 1:2 and 1:3 in our tests. As Table 2 shows, it is clear that the original model outperforms both alternatives. Compared to the 1:2 ratio, the original model has more explanatory power in every way (i.e. higher overall accuracy, higher threat accuracy, higher non-threat accuracy). By contrast, the 1:3 ratio model has a slightly lower overall accuracy, marginally higher non-threat accuracy, but devastatingly lower threat accuracy (17.5%) than our original model because the increasing imbalance (from 7:10 to 1:3) turned the automated machine-learning process to a 'greedy' mode. As a result, it is our conclusion that the original model uses the golden ratio with the most accurate categorizing capacity.

**Table 2** Comparison with alternative models

| Models | Overall accuracy % (N) | Threat accuracy (Sensitivity) | Non-threat accuracy (Specificity) |
|---|---|---|---|
| Rare event I | 72.4% | 32.8% | 82.3% |
| 1:2 Ratio | (368/508) | (66/201) | (302/367) |
| Rare event II | 78.8% | 17.5% | 99% |
| 1:3 Ratio | (656/832) | (36/206) | (620/626) |
| Out of sampling | 54.1% | 14.5% | 82.1% |
| : KJI → KJU | (647/1195) | (72/495) | (575/700) |
| NK leadership | KJI: 68.5% (98/143) | 45.9% (28/61) | 79.3% (65/82) |
| : KJI vs. KJU | KJU: 59.5% (195/328) | 61.2% (93/152) | 58% (102/176) |
| SK leadership | Con: 61.4% (221/360) | 36.3% (58/160) | 81.5% (163/200) |
| : Con. vs. Prog. | Prog.: 68.1% (98/144) | 49.3% (35/71) | 86.3% (63/73) |
| No cheonan | 66.9% | 51.4% | 78.7% |
| Case | (273/408) | 92/178 | 181/230 |
| Our Model | 82.3% | 50.3% | 96.1% |
| | (401/487) | (74/147) | (327/340) |

Second, it is also necessary to check whether an out-of-sample approach is a better alternative to our original model. In our analysis, we developed a model by using 70% of data from the five North Korean attacks and then applied it to the remaining 30% of data from. One may suspect, however, that a better approach is to elaborate a model from the first three North Korean strikes (i.e. during the Kim Jong-il period) and then apply it to the remaining two attacks (i.e. during the Kim Jong-un era) in order to see whether it yields more explanatory power. As Table 2 shows, the opposite turns out to be the case. In fact, the out-of-sampling model loses explanatory power in every possible way (i.e. lower overall accuracy, lower threat accuracy, and lower non-threat accuracy). As a result, our original model outperforms the out-of-sampling alternative.

Third, it is necessary to check whether or not there is a leadership change effect in North Korea. Has the change of power from Kim Jong-il to Kim Jong-un produced such significant differences in their leadership style that it may be worth developing two separate models (one for the Kim Jong-il period the other for the Kim Jong-un era) instead of a single model covering the entire period as we did? As

Table 2 shows, the double platform alternative yields a worse outcome overall. In fact, its only advantage is that the Kim-Jong-un-model has slightly higher threat accuracy (61.2%) than our original model (50.3%). In every other aspect, however, the Kim-Jong-un-model performs much worse. Moreover, the Kim-Jong-il-model has a much lower accuracy in all respects than our original model. As a result, it is clear that the double platform is a worse alternative to our single platform. The recent leadership change in Pyongyang has shown little effects as far as its military provocations are concerned.

Fourth, it is important to investigate whether a leadership change in South Korea has any impact on North Korean military provocations. Has the oscillation of power between the 'progressive' administrations (Kim Dae-jung and Roh Moo-hyun) and the 'conservative' regimes (Lee Myung-bak and Park Geun-hye) invoked different responses from Pyongyang? If so, we should expect a double platform (one model for the conservative period vs. the other model for the progressive era in South Korea) to outperform the original single platform. As Table 2 shows, however, the two separate models in the double-platform alternative perform much worse than the original single platform in all categories: an overall accuracy, threat accuracy, and non-threat accuracy. Clearly, the original model is a better choice. It is shown above that the leadership change in Pyongyang has not produced significant policy shifts in its military provocations. Likewise, leadership changes in South Korea do not seem to have invoked major policy shifts in Pyongyang as far as its military provocations are concerned.

Finally, it is worth examining an alternative model that excludes the sinking of the Cheonan naval ship case. Unlike the remaining four attacks in our dataset, the North Korean government has consistently denied its involvement in the Cheonan incident. If Pyongyang had not really sunk the Cheonan as it claimed, it means that our original model is performing below its full potential because it is built upon some erroneous cases (i.e. the Cheonan incident). In that case, we should expect the no-Cheonan alternative to outperform our original model. As Table 2 shows, however, when we exclude all the data related to the Cheonan case, the resulting model loses much explanatory power compared to the original model. Only in one category (a threat accuracy), the no-Cheonan-case alternative outperforms our original model, but even in that category, the difference is ignorable (50.3% vs. 51.4%). In
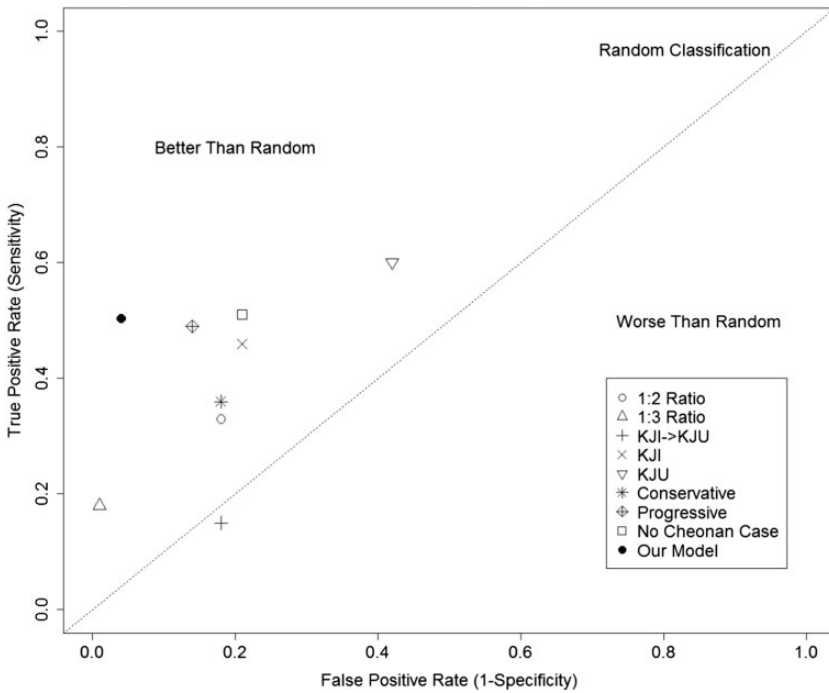
**Figure 6** Receiver operating characteristic (ROC) plot

all other categories, the original model shows much stronger performances: 82.3% vs. 66.9% in an overall accuracy and 96.1% vs. 78.7% in a non-threat accuracy. Although Pyongyang has denied its involvement in the Cheonan incident, the huge performance gap between the original model and the no-Cheonan alternative suggests otherwise.

A different way of comparing performances of alternative models is shown in Figure 6 where the Receiver Operating Characteristic (ROC) is presented. A ROC space is defined by a False Positive Rate (1 – Specificity) in the x-axis and a True Positive Rate (Sensitivity) in the y-axis. In a ROC space, a theoretically perfect model (i.e. a combination of 100% True Positive Rate with 0% False Positive Rate) appears in the upper left corner with the coordinates (0, 1). In comparison, a purely random model such as a coin toss is found along a 45-degree diagonal line. As a result, good models (i.e. those performing better than random guesses) are found above the 45 degree line (especially close to the y-axis), whereas bad models are located below it. In Figure 6, the

coordinates of our model (0.04, 0.5) means that it has a False Positive Rate of 4% while its True Positive Rate is 50%. By contrast, the coordinates of alternative models in the ROC space are as follows; (i) the Rare Event I model with 1:2 ratio is (0.18, 0.33); (ii) the Rare Event II model with 1:3 ratio is (0.01, 0.18); (iii) the Out of Sampling (KJI → KJU) model is (0.18, 0.15); (iv) the Leadership Change (KJI-only) model is (0.21, 0.46); (v) the Leadership Change (KJU-only) model is (0.42, 0.61); (vi) the model for Conservative South Korean leadership is (0.18, 0.36); (vii) the model for Conservative South Korean leadership is (0.14, 0.49); and (viii) the model that excludes the Cheonan case is (0.21, 0.51). As one can see, all eight alternatives yield their coordinates either close to or lower than the diagonal line in the ROC space, whereas our model lies above the diagonal line and close to the y-axis in the ROC space. As a result, we can conclude that the original model performs better than its potential rivals.

# 4 Conclusion

In this article, we studied articles published by the KCNA, the official news outlet of North Korea, in order to analyze patterns of its conventional military provocations. To this end, we have adopted a new method of automated text classification through supervised machine learning. Our model investigated the frequency of all terms appearing in KCNA articles immediately prior to five North Korean military attacks between 1997 and 2013. The frequency of these terms was then compared with the frequency of terms appearing in KCNA articles published during peacetime without military provocations. The comparison brought to light a number of key terms – 'attack words' so to speak – whose appearances spiked in the KCNA prior to North Korean attacks. Based on these terms, our model correctly identifies eight of 10 articles as signs of imminent attacks or as peacetime news pieces.

Specifically, our model found five pattern-detecting terms of North Korean military threats: 'years', 'signed', 'assembly', 'June', and 'Japanese.' For a proper analysis of their meaning, we went back to the articles in which these five attack words appeared in order to investigate the contexts in which they were used by the North Korean government. Our investigation shows that in the lead-up to an attack,

Pyongyang displayed a strong tendency to emphasize the legacy of its military struggle against 'Japanese' colonialism and its fight against the U.S. imperialism during the Korean War (which began in 'June' 1950), perhaps in an effort to heighten a domestic patriotic fervor at a time of impending crisis. In addition, immediately before Pyongyang embarks on hostile provocations, the KCNA tends to increase its reprinting of 'signed' commentaries from *Rodong Sinmun*, which typically criticizes the United States, South Korea, and Japan in harsh terms, probably reflecting its deteriorating relations with the outside world.

It seems that our methodology is promising to studies of international conflicts of various sorts. First, as shown in this article, a machine-learning technology is useful in terms of detecting patterns that distinguish threats of Pyongyang from its 'noises' or 'bluffing.' Second, for other authoritarian regimes where there is a tight government controls over the mass media, our approach can be used to detect patterns, symptoms, or clues to escalating crisis. Finally, even for democratic countries with a free media, our approach can be used on various occasions. For instance, a machine-learning technology can be utilized to analyze certain aspects of domestic politics, such as signaling or communication between policymakers and the public. To cite a concrete example, we can use a machine-learning technique to examine how an aggressive foreign policy like 'sabre rattling' affects public perceptions of fear or crisis in democracy. Also, a machine-learning approach can be used to analyze external interactions of a democratic regime with other countries. For instance, we can test if there are any significant correlations between presidential elections in the US and the level of North Korean threats, or between North Korean nuclear tests and varying public reactions in South Korea.

As a final note, it is *not* our contention that the model developed in this article based on an automated machine-learning technique has detected *intentional* signals which Pyongyang sends to the outside world immediately prior to its military attack. Instead, the key attack words we have identified should be seen as signs or patterns that the North Korean regime *unwittingly* displays when it is inching toward a military option. For two reasons, we doubt an intentional or signaling nature of the attack words in our model. First, if Pyongyang has indeed used the five attack words as a deliberate signal to the outside world that it is gearing up for military provocations, why would it send the signal in

such an arcane way that can be detected only through a complicated automated text classification process? If the North Korean regime intends to send a signal to the outside world, there is a better, clearer, and more visible option, such as an Official Announcement by the Ministry of Foreign Affairs, which Pyongyang has occasionally published in the pages of the KCNA on important issues. Second, if it had been sending intentional signals of an impending military strike, why would the North Korean regime deny such an attack afterwards? If repeated, an *ex post* denial would only reduce the credibility of an *ex ante* signal, creating an image of North Korea as a casual bluffer. For their arcane nature and occasional *ex post* denial, the five attack words in our model should not be treated as a premeditated signal from Pyongyang. Instead, they should be understood as inadvertent signs or patterns unconsciously displayed by the North Korean government prior to a military attack. In fact, it is the unplanned nature of these attack words that provides even more valuable information to the outside world because it eliminates the possibility of feigned or false signals from North Korea. Like a pitcher who unknowingly flinches before he throws a fast ball, Pyongyang may unwittingly display certain patterns before it launches a military strike.

## Acknowledgements

## References

Breslow, N.E. (1196) 'Statistics in epidemiology: the case-control study.' *Journal of American Statistical Association*, 91, 433, 14–28.

Cha, V. (2010) The Sinking of the Cheonan. *Center for Strategic & International Studies.* (2010, April 22). http://csis.org/publication/sinking-cheonan (24 July 2014, date last accessed).

Choe, S.-H. (2009) Korean Navies Skirmish in Disputed\Waters. *New York Times.* (2009, November 10). http://www.nytimes.com/2009/11/11/world/asia/11korea.html?_r=0 (23 June 2014, date last accessed).

Ertekin, S. (2013) 'Adaptive oversampling for imbalanced data classification.', in G. Erol & L. Ricard (ed.), *Information Sciences and Systems*, pp 261–269. New York, Springer.

Global Security. (2002) The naval clash on the yellow sea on 29 June 2002 between South and North Korea: the situation and ROK's position. (2002, July 1). *Global Security.* http://www.globalsecurity.org/wmd/library/news/dprk/2002/dprk-020701-1.htm (22 July 2014, date last accessed).

Hong, S. (2012) 'North Korea's capability to conduct provocations and ROK-US capability to counter them.' *Korea Association of Defense Industry Studies*, 19, 2, 135–136.

Hopkins, D. and King, G. (2010) 'Extracting systemic social science meaning from text.' *American Journal of Political Science*, 54, 1, 229–47.

Jung, S.-C. (2013) *Internal Instability and External Provocation.* Seoul: KINU.

Jung, S.-Y. (2008) 'A study on the origin of the pueblo incident on 1968.' *Kukjejŏngchiyŏn'gu*, 11, 2, 179–207.

Jurafsky, D. and Martin, J. (2009) *Speech and Natural Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition.* NJ: Prentice Hall.

Kang, C.-G. (2013) 'A laboratory for recursive party.' *Hangukcontentshakhoe*, 11, 4, 23–32.

King, G. and Zeng, L. (2001a) 'Logistic regression in rare events data.' *Political Analysis*, 9, 2, 137–163.

King, G. and Zeng, L. (2001b) 'Explaining rare events in International Relations.' *International Organization*, 55, 3, 693–715.

Ko, M.-K. (2015) 'North Korean military adventurism in the late 1960s and the changes of party-military relations.' *Hyŏndaibukhanyŏngu*, 18, 3, 7–58.

Lee, W.-K. (2014) 'A case study on the provocations by NK.' *Kunsaji*, 91, 6, 63–110.

Litwak, R. (2007) *Regime Change.* Washington, D.C.: Johns Hopkins University Press.

Macfie, N. (2013) The Battles of the Korean West Sea. (2010, November 29). *Reuters.*http://www.reuters.com/article/2010/11/29/us-korea-north-clashes-idUSTRE6AS1AL20101129 (5 August 2014, date last accessed).

Mearsheimer, J. J. (2001). *The Tragedy of Great Power Politics.* New York, WW Norton & Company.

Moore, M. and Hutchison, P. (2010) Yeonpyeong Island: A History.(2010, November 23). *The Telegraph.* http://www.telegraph.co.uk/news/worldnews/

asia/southkorea/8155486/Yeonpyeong-Island-A-history.html (7 August 2014, date last accessed).

Oh, I.-W. (2011) 'South Korea's countermeasures against North Korea's Armed Provocation and Offensive dialogue proposal.' *T'ongiljyŏllyak*, 11, 1, 227–266.

Rich, T. (2012a) 'Like father like son? Correlates of leaderhip in North Korea's english language news.' *Korea Observer*, 43, 4, 649–674.

Rich, T. (2012b) 'Deciphering North Korea's nuclear rhetoric: an automated content analysis of KCNA news.' *Asian Affairs: An American Review*, 39, 73–89.

Sigal, L. (1998) *Disarming Strangers: Nuclear Diplomacy with North Korea*. Princeton: Princeton University Press.

Sohn, J.-Y. (2002) South, North Korea clash at sea. *CNN* (2002, June 29). http://edition.cnn.com/2002/WORLD/asiapcf/east/06/29/korea.warships/. (21 July 2014, date last accessed).

Sudworth, J. (2010) How South Korean Ship was Sunk. *BBC News* (2010, May 20). http://www.bbc.co.uk/news/10130909 (4 August 2014, date last accessed).