

1
2
3
4
5
6
7
8
9
10
11
12
13

Text classification to streamline online wildlife trade analyses

Oliver C. Stringham^{1,2*}, Stephanie Moncayo¹, Katherine G.W. Hill¹, Adam Toomes¹, Lewis Mitchell²,
Joshua V. Ross², Phillip Cassey¹

1. Invasion Science & Wildlife Ecology Lab, University of Adelaide, SA 5005, Australia
2. School of Mathematical Sciences, University of Adelaide, SA 5005, Australia

*Corresponding Author: Oliver C. Stringham, The University of Adelaide, North Terrace Campus,
Adelaide, SA, 5000, Australia. oliverstringham@gmail.com

Running headline: Text classification and online wildlife trade

Keywords: data cleaning, internet, natural language processing, text mining, wildlife trade

14 **Abstract**

- 15 1. Automated monitoring of websites that trade wildlife is increasingly necessary to inform
16 conservation and biosecurity efforts. However, e-commerce and wildlife trading websites can
17 contain a vast number of advertisements, an unknown proportion of which may be irrelevant to
18 researchers and practitioners. Given that many of these advertisements have an unstructured text
19 format, automated identification of relevant listings has not traditionally been possible, nor
20 attempted. Other scientific disciplines have solved similar problems using machine learning and
21 natural language processing models, such as text classifiers.
- 22 2. Here, we test the ability of a suite of text classifiers to extract relevant advertisements from an
23 Australian classifieds website where people can post advertisements of their pet birds (n = 16.5k
24 advertisements). Furthermore, in an attempt to answer the question ‘how much data is required to
25 have an adequately performing model?’, we conducted a sensitivity analysis by simulating
26 decreases in sample sizes to measure the subsequent change in model performance.
- 27 3. We found that text classifiers can predict, with a high degree of accuracy, which listings are relevant
28 (ROC AUC ≥ 0.98 , F1 score ≥ 0.77). From our sensitivity analysis, we found that text classifiers
29 required a minimum sample size of 33% (c. 5.5k listings) to accurately identify relevant listings (for
30 our dataset), providing a reference point for future applications of this sort.
- 31 4. Our results suggest that text classification is a viable tool that can be applied to the online trade of
32 wildlife to reduce time dedicated to data cleaning. However, the success of text classifiers will vary
33 depending on the advertisements and websites, and will therefore be context dependent. Further
34 work to integrate other machine learning tools, such as image classification, may provide better
35 predictive abilities in the context of streamlining data processing for wildlife trade related online
36 data.

37 **Introduction**

38 The global wildlife trade is a major concern for biodiversity conservation and biosecurity enforcement
39 (Smith et al. 2009). Information on the composition and volume of species, and where they are traded,
40 is highly valuable for informing conservation research and practice (Scheffers et al. 2019). The Internet is
41 an emerging source of data on the wildlife trade (Siriwat and Nijman 2020; Jarić et al. 2020).

42 Researchers, NGOs, and government agencies monitor websites that trade wildlife to quantify various
43 aspects of the trade (e.g., Sung & Fong 2018). Data gathered from the Internet are typically not
44 immediately ready for analysis (i.e., they are ‘messy’) and must be cleaned or processed to identify the
45 desired attributes for subsequent analysis (Dobson et al. 2020). This is especially true for classifieds,
46 forums, and social media sites where human users type their advertisements into an open (or ‘free
47 form’) text box. Consequently, relevant attributes cannot be extracted automatically (i.e., through web
48 scraping or computer-based data manipulation) due to non-uniformity across users’ advertisements
49 (different species names, abbreviations, misspelling, etc.) (Stringham et al. 2020). Likewise, depending
50 on the website, many online listings (i.e., posts) may contain items or taxa that are irrelevant for a given
51 research context. For instance, in a pet reptile forum, one can find users trading tanks, food, or other
52 accessories, which may not be relevant to researchers exploring the trade of live reptiles (e.g.,
53 Stringham and Lockwood 2018). The most common method to extract online wildlife trade data is to
54 manually inspect each listing and record the desired attributes. Depending on how many listings are
55 collected, the data cleaning process could represent an enormous amount of time and effort for
56 researchers. Wildlife-related web data is notorious for its scale: for example, Xu et al. (2018) tracked
57 around 140k tweets from a two-week period relevant to ivory and pangolin trade.

58

59 Automated methods of data cleaning such as machine-learning techniques and Natural Language
60 Processing (NLP) tools have potential to streamline the processing of wildlife trade data derived from

61 the Internet (Di Minin et al. 2019). A useful but unexplored application is to predict and extract relevant
62 online listings based on their text, which could save time in manual data processing steps if many
63 irrelevant listings exist in the dataset. In particular, text classification models relate the words associated
64 with a particular label, such as ‘relevant’ or ‘irrelevant’, to predict the label of an unknown data point. A
65 well-known application of text classifiers is filtering spam emails (Guzella and Caminhas 2009). In this
66 context, a text classification model uses a training dataset of labelled emails (spam or not spam) and
67 trains a model to predict those labels based on their constituent words. The resulting model labels new
68 incoming emails as spam or not. In the context of wildlife-trade data derived from the Internet, text
69 classification models have the potential to identify relevant listings and remove irrelevant listings that
70 do not sell wildlife (i.e., fish tanks, bird cages, food) by using the words in the listings. If shown to be
71 effective, this could save researchers substantial time in the data cleaning process.

72

73 Here, we examine if text classification models can predict which Internet listings are relevant to wildlife
74 trade research (for our own specific research purposes; e.g., Toomes et al. 2020). Further, to assist
75 future implementation of such models, we sought to identify how much data is needed for a text-
76 classification model to perform adequately well. We collected advertisement listings from a popular
77 Australian classifieds website where people trade their pet birds and accessories (e.g., bird cages or bird
78 toys). Bird trade is largely unregulated in Australia (but see Woolnough et al. 2020) and is highly diverse
79 with a large number of both native and alien species; with potential conservation and biosecurity
80 consequences (Vall-Ilosera & Cassey 2017). We observed three major categories of advertisements that
81 were irrelevant to our research objectives: (i) ‘junk’ listings (not trading birds); (ii) wanted
82 advertisements (requesting a bird); and (iii) the sale of domestic poultry – e.g., gamebirds, waterfowl,
83 and pigeons (non-target wildlife taxa). We manually labelled around 16.5k listings and tested the
84 efficacy of three commonly used text classification models at determining which listings were relevant

85 versus irrelevant. Next, we systematically removed records from our dataset and recorded the change in
86 model performance. Our results imply that text classification can be an incredibly useful time-saver
87 when cleaning data on the wildlife trade, which is structurally (textually) similar to the data we explore
88 here.

89 **Materials and Methods**

90 *Data collection and data curation*

91 We collected data from a popular Australian classifieds website daily over the course of five months (5
92 July 2019 to 5 December 2019). All information collected from the website was publicly available. We
93 received ethics approval from The University of Adelaide (ethics number: H-2020-184) to collect this
94 data and have anonymized the name of the website as good ethical practice (Hinsely et al. 2016). On the
95 website, people can post advertisements (i.e., listings) of items/animals they are trading. From each
96 listing, we collected: (i) the title; (ii) text description; (iii) date; and, (iv) images (if provided). The title and
97 text description fields are open text boxes where the user can type whatever they desire up to a
98 character limit. We collected a total of 66,704 unique listings. Given the large number of unique listings
99 collected, and the substantial resources required to manually clean the data, we labelled a random
100 subset of around 25% of the listings ($n = 16,509$). This took approximately 103 hours to label (at an
101 average rate of 161 listings per hour). Four different authors were labelers (SM, KH, AT, OS), and we did
102 not overlap labelling, although this is preferred practice (e.g. see Sheng et al. 2008).

103 For each listing, we manually labelled the taxa (e.g., species) being traded based on the title, the text,
104 and the pictures provided in the listing. Some listings contained more than one species being traded. We
105 identified the listing to the most specific taxonomic rank as possible (species or subspecies), but
106 occasionally not enough information was provided and the listing was identified to genus, family, order,

107 or class (i.e., bird). We resolved taxa names and obtained upper-level taxonomy using the Global
108 Biodiversity Information Facility database (GBIF 2020). For each listing, we recorded if the user was
109 requesting a bird species (i.e., a wanted advertisement), except in the case of domestic poultry species
110 (see below). We labelled listings not trading a live bird as ‘junk’ (i.e., bird accessory such as cage or bird
111 food).

112 *Preparing text for text classification models*

113 We considered all text written by the user (title and text description) for our analyses. To prepare or
114 ‘clean’ the text for the NLP text classification models, we followed standard NLP text cleaning
115 procedures (Silge and Robinson 2017) and removed special characters (emojis, dollar signs, numbers,
116 etc.), removed all punctuation, converted text to lowercase, and removed all numbers. Next, we
117 removed all stop words found in the following lexicons: *SMART*, *snowball*, and *onix*. We did not remove
118 the stop words: “want”, “wants”, “wanting”, or “wanted”, so we could distinguish wanted
119 advertisements. We stemmed each word using the Snowball stemmer. For the text classifier models, we
120 tokenized the text to be unigrams (i.e., one word) and did not consider further n-grams. Text cleaning
121 was performed in the statistical software R version (R Core Team 2020) using the following packages:
122 stringr (Wickham 2019), dplyr (Wickham et al. 2020), tidytext (Silge & Robinson 2016), and corpus (Perry
123 2020).

124 To test the classification of irrelevant listings (see ‘*Text classification models*’ below), we applied three
125 separate labels for each listing. The first label was for ‘junk’ listings, where a live bird was not being
126 traded (e.g., bird cage). The second label was for ‘wanted’ listings where a user was requesting a bird
127 species and not selling one. The final label was for taxa that we considered non-target for our purposes
128 (i.e., farm, poultry, or domesticated species). We called this label ‘domestic poultry’ and applied it to
129 listings that were selling birds in the taxonomic orders of Anseriformes (waterfowl) and Galliformes

130 (gamebirds) or trading domestic pigeons (*Columba livia domestica*). For text classification models, we
131 removed listings categorized as more than one label (i.e., 'domestic poultry' and 'wanted'). Further, for
132 the 'wanted' label, we removed listings if eggs were being advertised, as we did not simultaneously
133 record if egg advertisements were also labeled as 'wanted'. This resulted in a sample size (number of
134 listings used for text classification models) of 16,475 for 'domestic poultry', 16,446 for 'junk', and 13,751
135 for 'wanted'.

136 *Text classification models*

137 To classify irrelevant listings, we used three common supervised text classifiers: Logistic Regression,
138 Multinomial Naive Bayes, and Random Forest. At a basic level, each classifier considers each word (i.e.,
139 gram) and their frequency as a covariate (i.e., 'feature') (Bird et al. 2009). However, each classifier varies
140 in the algorithm used to classify observed listings as relevant or not (Bird et al. 2009). For each classifier,
141 the order of the words in the listing was unaccounted, thus earning the name 'bag of words' classifier.
142 We ran each model for each of the three labels mentioned above. We used 10-fold cross validation to
143 train the model and evaluate predictions. We used the cross-validated macro-average of the following
144 metrics to evaluate the performance of each model: receiver operating characteristic (ROC) curve and
145 its area under the curve (ROC AUC), precision-recall curve and its area under the curve (PR ROC),
146 precision, recall, negative predictive value (NPV), specificity, and F1 score (see Appendix S1 for more
147 information evaluation metrics). We extracted the top features (e.g., covariates) for each model. Text
148 classification models were performed in Python using the sci-kit learn library (Pedregosa et al. 2011),
149 while plotting was conducted in R using ggplot2 (Wickham 2016).

150 *Sensitivity analysis: degradation of model performance with diminishing sample size*

151 To test the sensitivity of model performance to changes in sample size, we implemented the text
152 classification model with iteratively smaller sample sizes. We systematically decreased the sample size
153 of the training set by 500 records at a time, removing at most 15k records (c. 91% of entire dataset). We
154 repeated this for each label and used 10-fold cross validation. To account for the variability in model
155 performance due to cross-validation, we repeated the text classification model for 100 iterations, for
156 each sample size explored. We recorded 10-fold cross validation statistics across each fold and model
157 iteration (1,000 values in total for each sample size). For this sensitivity analysis, we only considered the
158 logistic regression classifier and used the F1 value to evaluate model performance. We recorded the
159 maximum training set sample size at which the F1 score was 99% of its maximum value (i.e., the F1
160 score without reducing sample size).

161 **Results**

162 We manually categorized 16,509 listings, of which 15.0% (n=2,473) were labeled as 'junk', 21.9%
163 (n=3,615) were labeled as 'domestic poultry', 4.8% (n=787) were labelled as 'wanted' advertisements,
164 and the remaining (c. 58%) were 'for sale' advertisements of relevant bird taxa.

165
166 The text classifiers performed extremely well for the 'domestic poultry' label (Figure 1; Appendix S2),
167 with a cross-validated average ROC AUC of >0.99, Precision-Recall AUC of ≥ 0.97 , and F1 score of >0.95
168 for all text classifiers (Figures 1-3). The text classifiers for the 'junk' label also performed very well, with
169 marginally lower metric values compared to 'domestic poultry' (Figure 1). Further, all other metrics
170 evaluated suggested that the text classification models performed very well for these two labels (Figures
171 1-3; Appendix S2; see Appendix S3 for confusion matrices). The text classification models for the
172 'wanted' advertisement label performed less well, however, the Logistic Regression and Random Forest

173 classifiers for this label performed moderately well and each was much better than chance with a ROC
174 AUC > 0.98, Precision-Recall AUC > 0.88, and F1 score > 0.77. Overall, the 'wanted' classifiers were not as
175 good at predicting positive outcomes (e.g., if a listing is 'wanted'), yet did not struggle with predicting
176 negative outcomes (Specificity = 0.99, and Negative Predictive Value = 0.99 for Logistic Regression
177 classifier). In terms of relative performance between the classifiers, the Logistic Regression and Random
178 Forest classifiers slightly outperformed the Naive Bayes Classifier; however, overall, their performances
179 were comparable (Figure 1-3).

180
181 The top features for each label aligned with what should be expected and were similar across all text
182 classifiers (Figure 4). For the 'junk' label, grams such as "condit" (i.e., condition), "cage", "birdcag" (i.e.,
183 birdcage) were the top features. For the 'domestic poultry' label, grams such as "pigeon", "rooster", and
184 "chicken" were the top features. Finally, for the 'wanted' label, grams such as "want", "buy", "wtb" (an
185 acronym for 'want to buy'), and "unwant" (i.e. "unwanted") were the top features.

186
187 As we reduced the sample size of the training set, we observed a non-linear decrease in model
188 performance, where the F1 score initially declined gradually and then at an increasing rate at lower
189 sample sizes (Figure 5). There were differences in this decline in performance among labels. The
190 classifier for the 'domestic poultry' label realized 99% of the full model F1 score at c. 4.8k records (29%
191 of dataset). For the 'junk' label, this was c. 9.3k records and c. 6.3k records for the 'wanted' label (57%
192 and 45%, respectively). Stated another way, for the 'domestic poultry' label, the addition of c. 11k
193 labelled records from our manual data labelling only increased the model F1 score by 0.01. For 'junk'
194 and 'wanted', this value was c. 7.1k listings and c. 7.5k listings, respectively.

195

196 **Discussion**

197 Text classification can be a highly accurate method to extract relevant listings of wildlife found on the
198 Internet. In particular, for listings trading non-target taxa and listings trading bird accessories (e.g., bird
199 cages), text classification models were able to classify these listings with a very high degree of accuracy.
200 Although the performance of the model varied between labels, our results suggest that this technique
201 can be used to substantially lower the number of wildlife listings needed to be manually inspected, thus
202 saving considerable time and resources. Further, we provide clarity around the question of ‘how much
203 data is needed to guarantee an adequately performing model?’. Of the more than 16k listings we
204 manually labelled, our results suggest that, at most, only 9k listings were needed, although this number
205 varied by label.

206
207 Text classification models are commonplace in other disciplines and industries, which work heavily with
208 text data (e.g., Guzella & Caminhas 2009), yet have not been applied to data collected on the wildlife
209 traded occurring on the Internet. Importantly, from our dataset, around 60% of the listings were
210 relevant (for our purposes), representing a substantial amount of time and effort that would otherwise
211 be spent on manually removing irrelevant data. For the website we explored, we showed that text
212 classifiers predicted with great accuracy the advertisements that were not selling wildlife or selling non-
213 target wildlife. In particular, text classification models performed the best for identifying listings trading
214 non-target taxa (e.g., farm and domestic bird species). This is promising as a time saving tool because
215 sometimes the most commonly traded taxa are the ones of least interest to researchers (i.e., pigeons
216 and chickens; in our example). In contrast, the text classifiers had more difficulty distinguishing ‘wanted’
217 advertisements (where a user was requesting a bird) yet was still much better than chance. This
218 suggests that the words people used in wanted advertisements have some overlap with those words
219 used in non-wanted advertisements (e.g., names of species), and thus yields lower predictive abilities.

220 Importantly, we demonstrate that the model performance will likely not improve with more data
221 because we observed a plateau of model performance after around 6k listings for wanted
222 advertisements (45% of sample size). Therefore, even if we manually labelled many more listings, the
223 model performance is unlikely to increase. This highlights an important point that model performance is
224 a function of the underlying data itself (i.e., text) and not of the lack of data (once an adequate sample
225 size is achieved).

226

227 How much data is required for an adequately performing text classification model? Our results show
228 that this number will vary by what is being classified. For this study, we cleaned a substantial number of
229 listings (c. 16k) yet found that model performance marginally increased after 5k to 9k records (31% to
230 56% of total effort). Thus, for other researchers who may not have the resources to invest this much
231 effort, or are looking for a more efficient way to curate messy online data, our results provide guidance
232 on how much data is needed before text classification can be used. We recommend establishing
233 computer code to test the model performance and then repeatedly check the model performance at
234 regular intervals (e.g., every 1k records cleaned). Ultimately, the labelled dataset will need to
235 encapsulate the variation of words (i.e., vocabulary) used for a particular label for the text classifier to
236 perform well. For instance, for the 'junk' label, the model performance plateaued at around 5k more
237 records than it did for other labels. We hypothesize the words that Internet users write for the listings
238 that fall under the 'junk' label has more variation (i.e., more words) and thus, we needed a larger sample
239 size of labelled listings to account for that variation.

240

241 An important limitation of text classification (and other machine learning tools) is that they are highly
242 context dependent (Lambda et al. 2019). Our specific classifiers were developed based on the text of
243 birds being traded online in Australia and will likely be less useful for birds being traded in other

244 countries and almost entirely useless if looking at other taxa (e.g., fish or plants) or in another language.
245 The reason for this lack of generalization is because words used, and their frequency, will vary under
246 different contexts. For instance, when looking at the trade of aquarium fish, a common irrelevant
247 advertisement may be the sale of a fish tank, something that is not found when trading birds. We
248 recommend that researchers consider each context separately when using these tools. Since manual
249 data processing is likely always required to analyze the data, these tools can be tested throughout the
250 cleaning stage to see if applicable.

251
252 Besides extracting relevant advertisements, text classifiers have the potential to identify the species
253 being traded in online advertisements. Our results suggest that this will be possible for commonly
254 traded taxa, with large amounts of data. For instance, in our study, advertisements for a group of
255 species (waterfowl, gamebirds, and pigeons) comprised around 3.6k listings (22% of dataset) and were
256 highly distinguishable using the text classifiers. The same kinds of models can be used to identify
257 individual species of interest; however, text classifiers (like all machine learning techniques) require a
258 large volume of data to perform well (Di Minin et al. 2019). In many cases, individual species of interest
259 may not have enough advertisements to build adequate text classifiers. Thus, alternative methods such
260 as matching species names (scientific, common, or trade names) to the text of advertisements using a
261 fuzzy string-matching model (e.g., Levenshtein distance) may yield better results. In fact, if consistent
262 patterns are used by users (e.g., the same species name is used by many users), string matching may
263 yield just as good or better results than text classifiers. While our study relied exclusively on the text of
264 the advertisement, there are other attributes of an Internet listing that can be considered for automated
265 cleaning. For instance, a related study used metadata attributes of online listings (e.g., the number of
266 views and the price) to classify illegal sales of elephant ivory (Hernandez-Castro & Roberts 2015). In
267 cases with no or limited text provided (e.g., only a photo is posted), machine learning techniques such as

268 image classification could assist in the classification of species or the product traded (Norouzzadeh et al.
269 2018). Integrating text classification with the aforementioned models may improve predictive ability,
270 and we recommend this as a future area of research and development for the wildlife trade related
271 online data.

272

273 Given that a substantial proportion of online listings may not be relevant to wildlife trade research (e.g.,
274 40% irrelevant for our dataset), text classification methods can substantially decrease the amount of
275 time spent processing raw data. Here, we demonstrate that text classification can be viable tool to
276 identify irrelevant listings. When considering data on the scale of ‘big data’ of tens to hundreds of
277 thousands of online advertisements (e.g. Olden et al 2020), text classifiers have the potential to save
278 tens to hundreds of hours of curation effort. We recommend future application of text classifiers and
279 testing other machine learning and natural language processing tools when cleaning messy data
280 collected from the Internet on wildlife trade.

281 **Data Availability**

282 Data and code for text classification are available from the *figshare* repository at

283 <https://doi.org/10.6084/m9.figshare.14032742> and from GitHub at

284 https://github.com/ocstringham/text_classification_wildlife_trade/.

285

286 **Acknowledgements**

287 This research was funded by the Centre for Invasive Species Solutions (Project PO1-I-002).

288 **Author Contributions**

289 OCS, LM, JVR, and PC conceived the ideas and designed the methodology. OCS collected the data. SM,

290 KGWT, and AT cleaned the data. OCS analyzed the data and led the writing of the manuscript. All

291 authors contributed critically to the drafts and gave final approval for publication.

292 **References**

- 293 Bird, S., Klein, E., & Loper, E. (2009). *Natural Language Processing with Python: Analyzing Text with the*
294 *Natural Language Toolkit*. O'Reilly Media, Inc.
- 295 Dobson, A. D. M., Milner-Gulland, E. J., Aebischer, N. J., Beale, C. M., Brozovic, R., Coals, P., Critchlow, R.,
296 Dancer, A., Greve, M., Hinsley, A., Ibbett, H., Johnston, A., Kuiper, T., Le Comber, S., Mahood, S. P.,
297 Moore, J. F., Nilsen, E. B., Pocock, M. J. O., Quinn, A., ... Keane, A. (2020). Making Messy Data Work
298 for Conservation. *One Earth*, 2(5), 455–465. <https://doi.org/10.1016/j.oneear.2020.04.012>
- 299 Guzella, T. S., & Caminhas, W. M. (2009). A review of machine learning approaches to Spam filtering.
300 *Expert Systems with Applications*, 36(7), 10206–10222.
301 <https://doi.org/10.1016/j.eswa.2009.02.037>
- 302 Hernandez-Castro, J., & Roberts, D. L. (2015). Automatic detection of potentially illegal online sales of
303 elephant ivory via data mining. *PeerJ Computer Science*, 1, e10. [https://doi.org/10.7717/peerj-](https://doi.org/10.7717/peerj-cs.10)
304 [cs.10](https://doi.org/10.7717/peerj-cs.10)
- 305 Hinsley, A., Lee, T. E., Harrison, J. R., & Roberts, D. L. (2016). Estimating the extent and structure of trade
306 in horticultural orchids via social media. *Conservation Biology*, 30(5), 1038–1047.
307 <https://doi.org/10.1111/cobi.12721>
- 308 Jarić, I., Correia, R. A., Brook, B. W., Buettel, J. C., Courchamp, F., Di Minin, E., Firth, J. A., Gaston, K. J.,
309 Jepson, P., Kalinkat, G., Ladle, R., Soriano-Redondo, A., Souza, A. T., & Roll, U. (2020). iEcology:
310 Harnessing Large Online Resources to Generate Ecological Insights. *Trends in Ecology & Evolution*,
311 35(7), 630–639. <https://doi.org/10.1016/j.tree.2020.03.003>
- 312 Lamba, A., Cassey, P., Segaran, R. R., & Koh, L. P. (2019). Deep learning for environmental conservation.
313 *Current Biology*, 29(19), R977–R982. <https://doi.org/10.1016/j.cub.2019.08.016>

314 Minin, E. D., Fink, C., Hiippala, T., & Tenkanen, H. (2019). A framework for investigating illegal wildlife
315 trade on social media with machine learning. *Conservation Biology*, 33(1), 210–213.
316 <https://doi.org/10.1111/cobi.13104>

317 Norouzzadeh, M. S., Nguyen, A., Kosmala, M., Swanson, A., Palmer, M. S., Packer, C., & Clune, J. (2018).
318 Automatically identifying, counting, and describing wild animals in camera-trap images with deep
319 learning. *Proceedings of the National Academy of Sciences*, 115(25), E5716–E5725.
320 <https://doi.org/10.1073/pnas.1719367115>

321 Nunez-Mir, G. C., Iannone, B. V., Pijanowski, B. C., Kong, N., & Fei, S. (2016). Automated content
322 analysis: Addressing the big literature challenge in ecology and evolution. *Methods in Ecology and*
323 *Evolution*, 7(11), 1262–1272. <https://doi.org/10.1111/2041-210X.12602>

324 Olden, J. D., Whattam, E., & Wood, S. A. (2020). Online auction marketplaces as a global pathway for
325 aquatic invasive species. *Hydrobiologia*. <https://doi.org/10.1007/s10750-020-04407-7>

326 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer,
327 P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., &
328 Duchesnay, É. (2011). Scikit-learn: Machine Learning in Python. *Journal of Machine Learning*
329 *Research*, 12(85), 2825–2830.

330 Perry P.O. (2020). corpus: Text Corpus Analysis. R package version 0.10.1. [https://CRAN.R-](https://CRAN.R-project.org/package=corpus)
331 [project.org/package=corpus](https://CRAN.R-project.org/package=corpus)

332 R Core Team (2020). R: A language and environment for statistical computing. R Foundation for
333 Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

334 Scheffers, B. R., Oliveira, B. F., Lamb, I., & Edwards, D. P. (2019). Global wildlife trade across the tree of
335 life. *Science*, 366(6461), 71–76. <https://doi.org/10.1126/science.aav5327>

336 Sheng, V. S., Provost, F., & Ipeirotis, P. G. (2008). Get another label? Improving data quality and data
337 mining using multiple, noisy labelers. *Proceedings of the 14th ACM SIGKDD International*

338 *Conference on Knowledge Discovery and Data Mining*, 614–622.

339 <https://doi.org/10.1145/1401890.1401965>

340 Silge, J., & Robinson, D. (2017). *Text Mining with R: A Tidy Approach* (1st edition). O'Reilly Media.

341 Siritwat, P., & Nijman, V. (2020). Wildlife trade shifts from brick-and-mortar markets to virtual
342 marketplaces: A case study of birds of prey trade in Thailand. *Journal of Asia-Pacific Biodiversity*.
343 <https://doi.org/10.1016/j.japb.2020.03.012>

344 Smith, K. F., Behrens, M., Schloegel, L. M., Marano, N., Burgiel, S., & Daszak, P. (2009). Reducing the
345 Risks of the Wildlife Trade. *Science*, 324(5927), 594–595. <https://doi.org/10.1126/science.1174460>

346 Stringham, O. C., & Lockwood, J. L. (2018). Pet problems: Biological and economic factors that influence
347 the release of alien reptiles and amphibians by pet owners. *Journal of Applied Ecology*, 55(6),
348 2632–2640. <https://doi.org/10.1111/1365-2664.13237>

349 Stringham, O. C., Toomes, A., Kanishka, A. M., Mitchell, L., Heinrich, S., Ross, J. V., & Cassey, P. (2020). A
350 guide to using the Internet to monitor and quantify the wildlife trade. *Conservation Biology*.
351 <https://doi.org/10.1111/cobi.13675>

352 Sung, Y.-H., & Fong, J. J. (2018). Assessing consumer trends and illegal activity by monitoring the online
353 wildlife trade. *Biological Conservation*, 227, 219–225.
354 <https://doi.org/10.1016/j.biocon.2018.09.025>

355 Toivonen, T., Heikinheimo, V., Fink, C., Hausmann, A., Hiippala, T., Järvi, O., Tenkanen, H., & Di Minin, E.
356 (2019). Social media data for conservation science: A methodological overview. *Biological
357 Conservation*, 233, 298–315. <https://doi.org/10.1016/j.biocon.2019.01.023>

358 Toomes, A., Stringham, O. C., Mitchell, L., Ross, J. V., & Cassey, P. (2020). Australia's wish list of exotic
359 pets: Biosecurity and conservation implications of desired alien and illegal pet species. *Neobiota*,
360 60, 43–59. <https://doi.org/10.3897/neobiota.60.51431>

361 Vall-Ilosera, M., & Cassey, P. (2017). Leaky doors: Private captivity as a prominent source of bird
362 introductions in Australia. *PLOS ONE*, 12(2), e0172851.
363 <https://doi.org/10.1371/journal.pone.0172851>

364 Wickham H. (2016) *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.

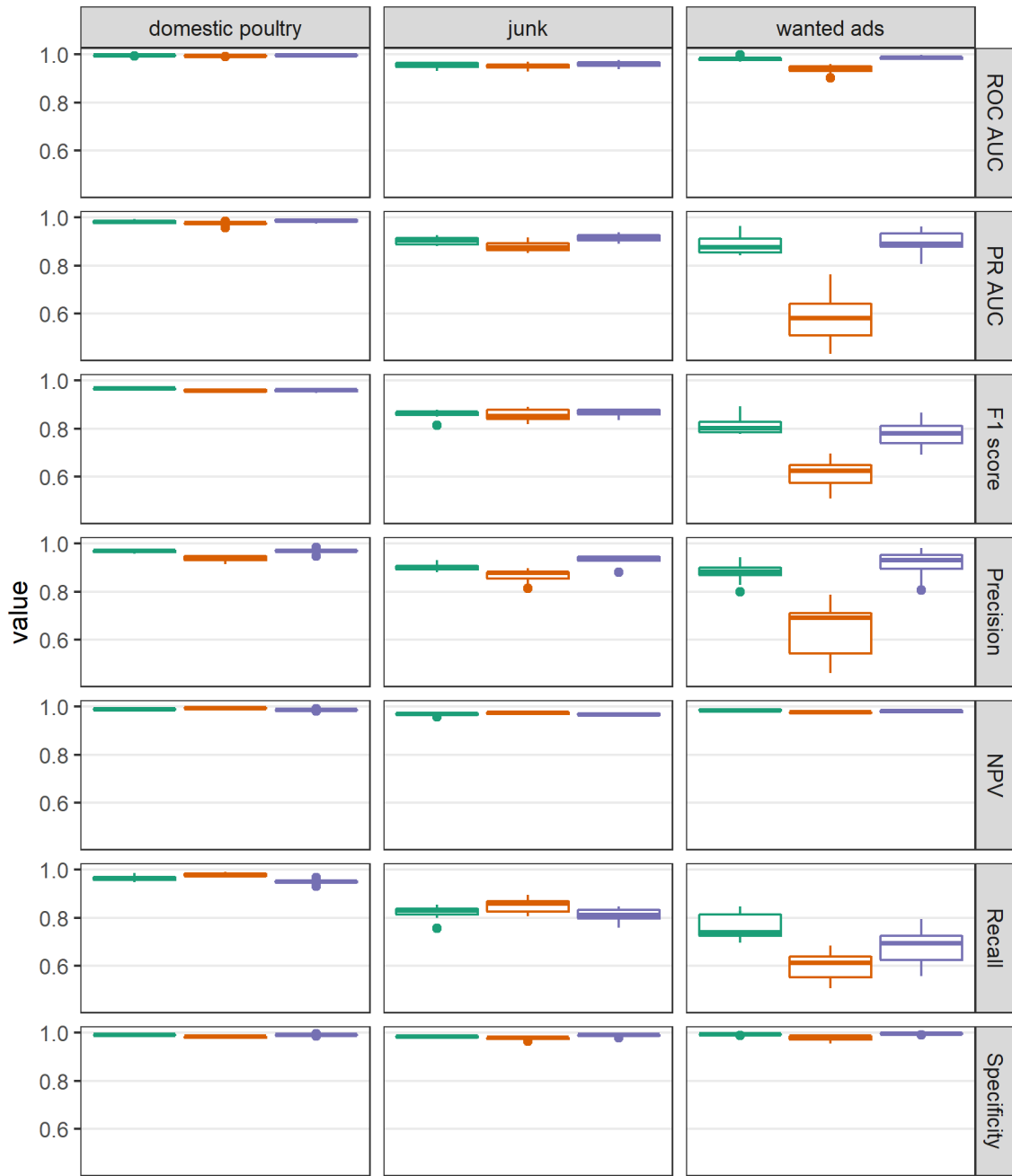
365 Wickham H. (2019). *stringr: Simple, Consistent Wrappers for Common String Operations*. R package
366 version 1.4.0. <https://CRAN.R-project.org/package=stringr>

367 Wickham H., François R., Henry L., Müller K. (2020). *dplyr: A Grammar of Data Manipulation*. R package
368 version 1.0.0. <https://CRAN.R-project.org/package=dplyr>

369 Woolnough, A.P., de Milliano J.M., van Polen Petel T., and Cassey P. (2020). A policy approach to
370 nonindigenous bird management in Victoria – managing potential threats to biodiversity, social
371 amenity and economic values. *Victorian Naturalist*. 137 (6); 203-209

372 Xu, Q., Li, J., Cai, M., & Mackey, T. K. (2019). Use of Machine Learning to Detect Wildlife Product
373 Promotion and Sales on Twitter. *Frontiers in Big Data*, 2. <https://doi.org/10.3389/fdata.2019.00028>

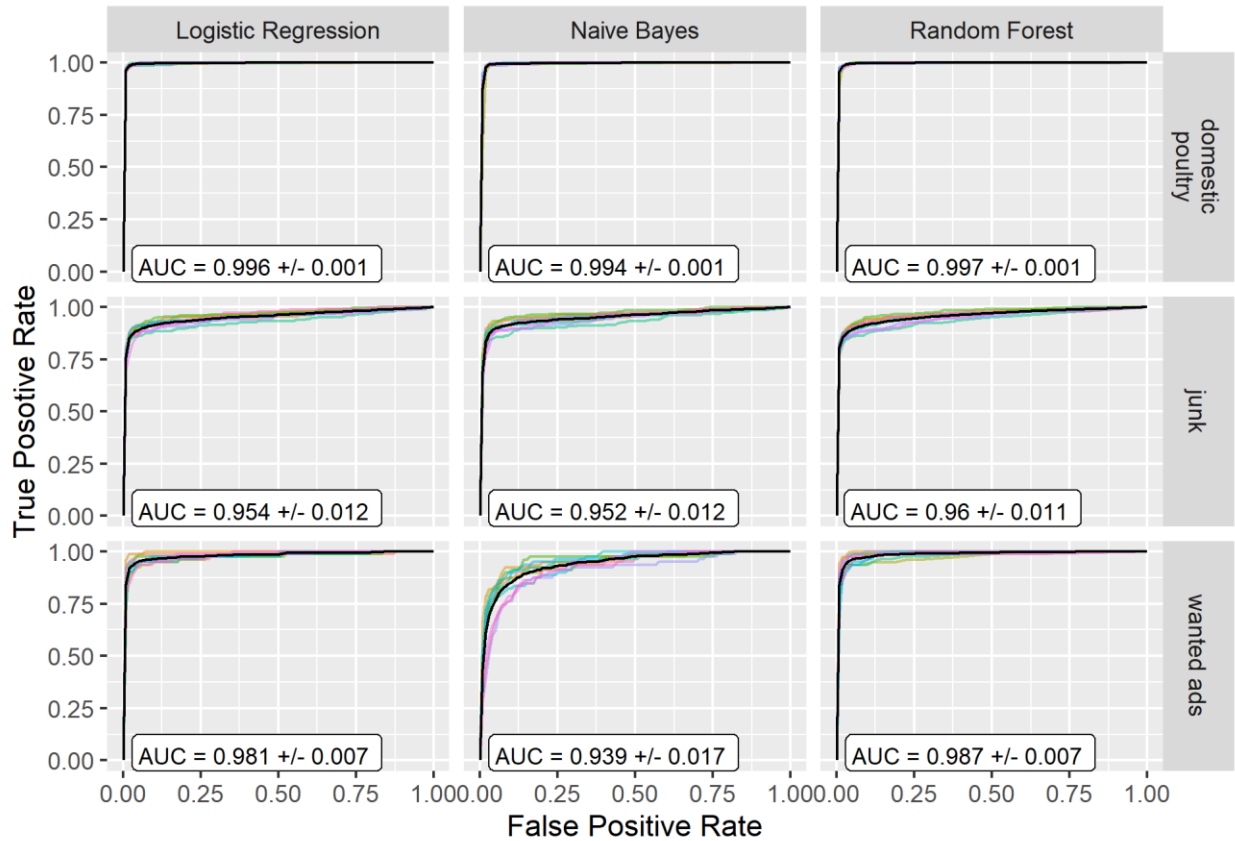
374



Classifier  Logistic Regression  Naive Bayes  Random Forest

377 *Figure 1.*

378 Model evaluation metrics (rows) across 10 cross-validation folds using different text classifiers evaluated
379 for three different labels (columns). See Appendix S1 for more information and calculation of the
380 evaluation metrics and Appendix S2 for exact metric values.



381

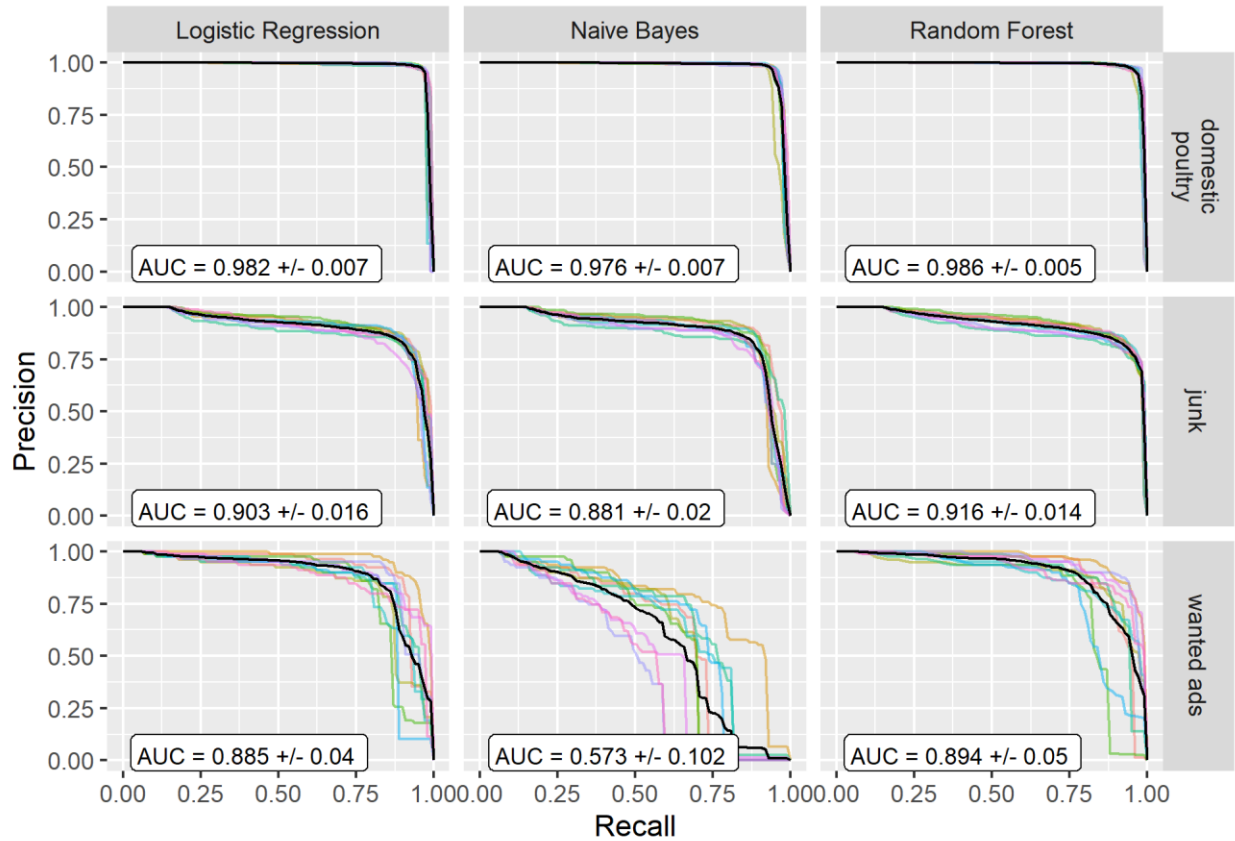
382 *Figure 2.*

383 Receiver operating characteristic curves and the area under the curve (ROC AUC). Three different text

384 classifiers (columns) were tested across three different labels (rows). For each panel, each line

385 represents one cross-validation fold and the solid black line represents the average across all cross-

386 validation folds. Average AUC (area under curve) values are reported with standard deviation.



387

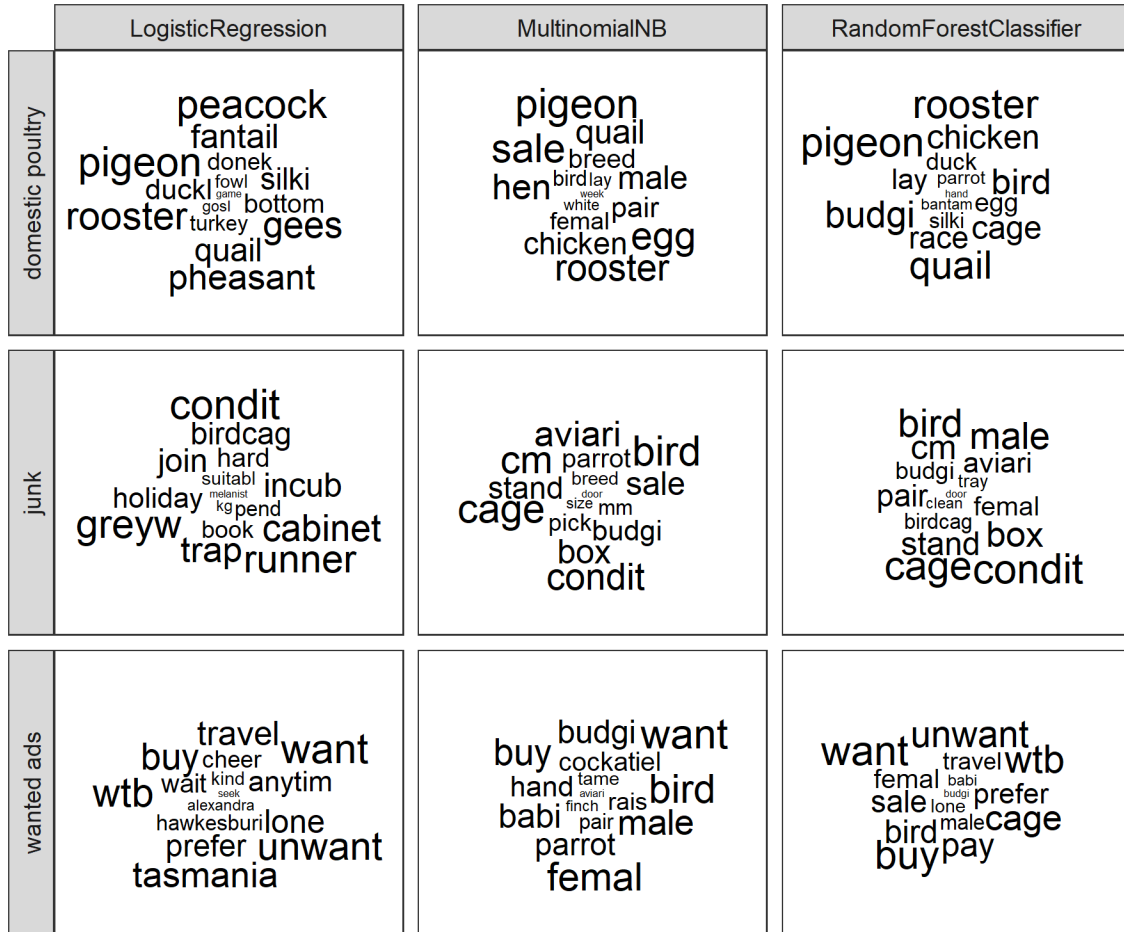
388 *Figure 3.*

389 Precision recall curves and the area under the curve (PR AUC). Three different text classifiers (columns)

390 were tested across three different labels (rows). For each panel, each line represents one cross-

391 validation fold and the solid black line represents the average across all cross-validation folds. Average

392 AUC (area under curve) values are reported with standard deviation.



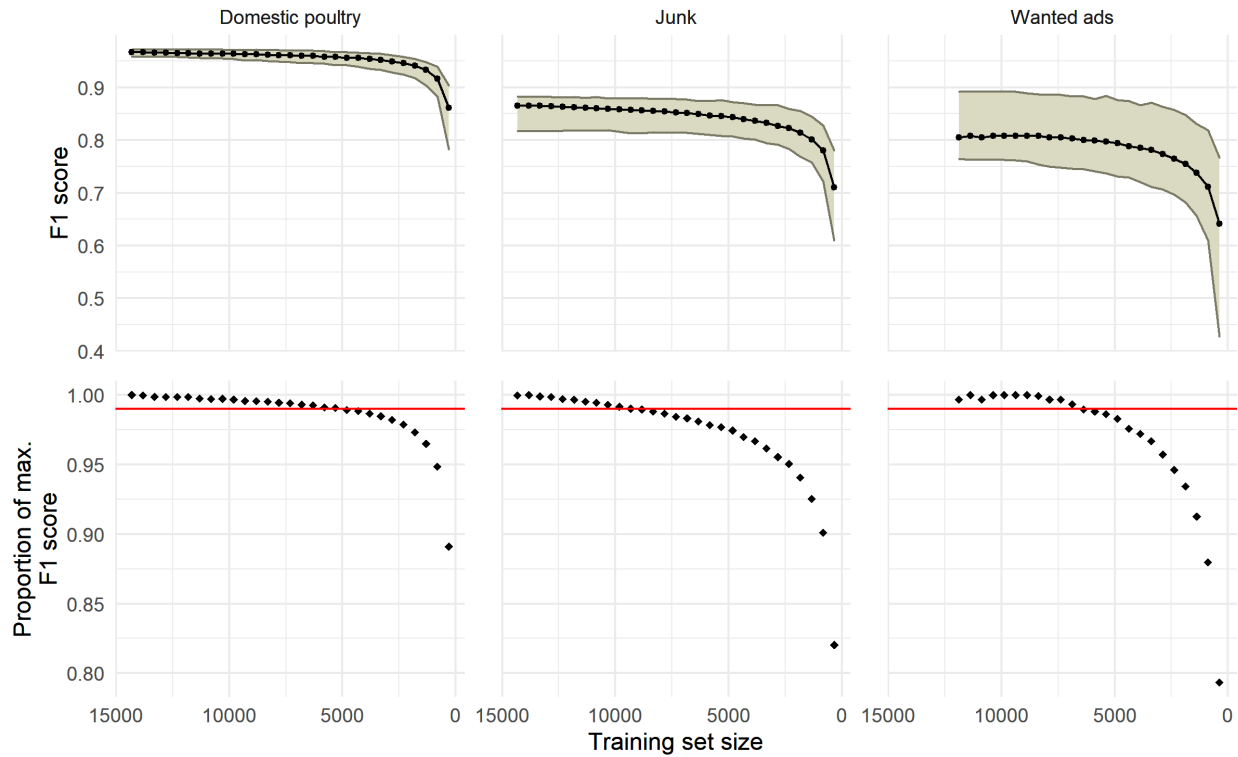
393

394 *Figure 4.*

395 Word clouds of top words (i.e., features or grams) for each label (rows) and classifier (columns). The size

396 of the word corresponds to importance, where larger words indicate higher importance. Note that

397 words are stemmed (e.g., condition is stemmed to condit).



398

399 *Figure 5.*

400 The effects of reducing sample size on text-classifier model performance. Top row: the F1 score
 401 evaluated at decreasing sample size (training set) values. Ribbons represent the 95% quantile range
 402 from 100 iterations of 10-fold cross validation logistic regression text classification, repeated for each
 403 specified label ('domestic poultry', 'junk', and 'wanted'). Bottom row: the proportion of the maximum
 404 F1 score, evaluated at each sample size, for each label. Only the median value was considered. The red
 405 horizontal line represents 0.99 of the maximum F1 score.

406 **Supporting Information**

- 407 • Appendix S1: Definitions of metrics used
- 408 • Appendix S2: Table of model metrics
- 409 • Appendix S3: Confusion matrices for all models

410 *Appendix S1: Definitions of metrics used*

411 Confusion matrix derived metrics

412 We evaluated several commonly used machine learning diagnostic metrics derived from confusion
413 matrix values (Appendix S3): true positives (TP), false negatives (FN), true negatives (TN), and false
414 positives (FP) (Fielding and Bell 1997). *Precision* is the proportion of correctly predicted positives
415 compared to all predicted positives. *Recall* is the proportion of correctly predicted positives compared to
416 all observed positives. The *Negative predictive value* is the proportion of correctly predicted negatives
417 compared to all predicted negatives. Finally, the *Specificity* is the proportion of correctly predicted
418 negatives compared to all observed negatives. Mathematically, the metrics are defined (Fielding and
419 Bell 1997) as follows:

420
$$Precision = \frac{TP}{TP + FP}$$

421
$$Recall = \frac{TP}{TP + FN}$$

422
$$Negative\ predictive\ value = \frac{TN}{FN + TN}$$

423
$$Specificity = \frac{TN}{FP + TN}$$

424 Further, the *F1 score*, is defined as the harmonic mean of precision and recall, mathematically:

425
$$F1\ score = 2 \cdot \frac{Precision \cdot Recall}{Precision + Recall}$$

426

427 Receiver operating characteristic (ROC) curve

428 The ROC curve shows the performance of a classification model at varying classification thresholds
429 (Fewcett 2006). The curve plots two metrics: False positive rate (i.e., $1 - Specificity$) and True positive
430 rate (i.e., *Recall*). For each classification threshold (e.g., from 0.01 to 1.0 by units of 0.01), the false
431 positive rate and true positive rate are plotted (e.g., main text Figure 2). The area under the curve for

432 the ROC curve (ROC AUC) is a measure of the positive predictive ability of the classification model (e.g.,
433 the ability to predict true positives versus false positives), where an ROC AUC of 0.5 represents positive
434 predictive ability equivalent to chance and an ROC AUC of 1 represents perfect positive predictive
435 ability.

436

437 Precision-Recall (PR) Curve

438 Like the ROC curve, the precision-recall (PR) curve also displays the performance of a classification
439 model at varying classification thresholds. However, for the PR curve, the tradeoff between *Precision*
440 and *Recall* is examined (not the True versus False positive rate examined in ROC curves). The PR curve is
441 useful when there are imbalanced class sizes (i.e., far fewer positives than negatives) because it does
442 not consider true positives in its calculation (Sofaer et al. 2018).

443

444 References

445 Fawcett, T. (2006). An introduction to ROC analysis. *Pattern Recognition Letters*, 27(8), 861–874.

446 <https://doi.org/10.1016/j.patrec.2005.10.010>

447 Fielding, A. H., & Bell, J. F. (1997). A review of methods for the assessment of prediction errors in
448 conservation presence/absence models. *Environmental Conservation*, 24(1), 38–49.

449 <https://doi.org/10.1017/S0376892997000088>

450 Sofaer, H. R., Hoeting, J. A., & Jarnevich, C. S. (2019). The area under the precision-recall curve as a
451 performance metric for rare binary events. *Methods in Ecology and Evolution*, 10(4), 565–577.

452 <https://doi.org/10.1111/2041-210X.13140>

453 *Appendix S2: Table of model metrics*

454 The macro-averaged values (10-fold cross validated) of model performance metrics for each label-classifier combination.

455

Classifier	ROC AUC	PR AUC	F1 score	Precision	Recall	NPV	Specificity
domestic poultry							
Logistic Regression	0.996	0.982	0.966	0.969	0.964	0.990	0.991
Naive Bayes	0.994	0.975	0.958	0.938	0.979	0.994	0.982
Random Forest	0.997	0.986	0.959	0.969	0.950	0.986	0.991
junk							
Logistic Regression	0.954	0.903	0.860	0.902	0.822	0.969	0.984
Naive Bayes	0.952	0.879	0.857	0.866	0.849	0.973	0.977
Random Forest	0.960	0.914	0.866	0.931	0.810	0.967	0.989
wanted							
Logistic Regression	0.981	0.886	0.815	0.878	0.764	0.986	0.993
Naive Bayes	0.939	0.579	0.614	0.641	0.600	0.976	0.978
Random Forest	0.987	0.893	0.775	0.913	0.676	0.981	0.996

456

457 *Appendix S3: Confusion “matrices”*

458 The median number (10-fold cross validated) of true positives, false negatives, false positives, and true negatives for each label-classifier
 459 combination.

Label	Classifier	True positive	False negative	False positive	True negative
domestic poultry	Logistic Regression	348	13	11	1272
	Naive Bayes	354	8	22	1262
	Random Forest	344	18	11	1272
junk	Logistic Regression	205	42	22	1378
	Naive Bayes	212	34	29	1371
	Random Forest	200	46	14	1386
wanted	Logistic Regression	58	20	8	1288
	Naive Bayes	48	30	22	1274
	Random Forest	54	24	4	1292

460