

1 Mapping Wildlife Species Distribution With Social 2 Media: Augmenting Text Classification With 3 Species Names

4 **Shelan S. Jeawak**¹

5 Cardiff University, School of Computer Science and Informatics, Cardiff, UK
6 JeawakSS@cardiff.ac.uk

7 **Christopher B. Jones**

8 Cardiff University, School of Computer Science and Informatics, Cardiff, UK
9 JonesCB2@cardiff.ac.uk

10 **Steven Schockaert**²

11 Cardiff University, School of Computer Science and Informatics, Cardiff, UK
12 SchockaertS1@cardiff.ac.uk

13 — Abstract —

14 Social media has considerable potential as a source of passive citizen science observations of the
15 natural environment, including wildlife monitoring. Here we compare and combine two main
16 strategies for using social media postings to predict species distributions: (i) identifying postings
17 that explicitly mention the target species name and (ii) using a text classifier that exploits all
18 tags to construct a model of the locations where the species occurs. We find that the first
19 strategy has high precision but suffers from low recall, with the second strategy achieving a
20 better overall performance. We furthermore show that even better performance is achieved with
21 a meta classifier that combines data on the presence or absence of species name tags with the
22 predictions from the text classifier.

23 **2012 ACM Subject Classification** I.2.6 Learning; I.2.1 Applications and Expert Systems

24 **Keywords and phrases** Social media, Text mining, Volunteered Geographic Information, Ecology

25 **Digital Object Identifier** 10.4230/LIPICs.GIScience.2018.<45>

26 **1** Introduction

27 The value of social media to assist in mapping and predicting geospatial phenomena has been
28 demonstrated in areas including the occurrence of disease, social unrest, natural disasters,
29 levels of wellbeing and characteristics of the man-made and natural environment [7, 8].
30 In the fields of environmental monitoring and wildlife observation there is clearly strong
31 potential for exploiting social media, reflected in the fact that searching for named species on
32 photo-sharing websites such as Flickr often reveals thousands of results, many of which are
33 associated with coordinates and almost all with time stamps. It can be envisaged that these
34 observations could complement the many effective citizen science campaigns that record
35 aspects of the natural environment and assist environmental scientists in understanding the
36 occurrence and behaviour of animals and plants [4]. Although many mentions of species
37 names in social media might not correspond to records of actual occurrences, several studies
38 have confirmed the validity of significant numbers of species observations in social media

¹ [has been sponsored by HCED Iraq.]

² [has been supported by ERC Starting Grant 637277.]



<45>:2 Mapping Wildlife Species Distribution

39 [1, 2]. While these studies highlight the potential value of such data, little progress has been
40 made to date on developing reliable automated methods for exploiting all the textual content
41 of social media postings for tasks such as mapping species distributions.

42 Here we present the results of experiments to predict species distribution based on
43 geocoded social media postings from the Flickr website. As a baseline approach we study
44 the performance of a method that predicts the occurrence of a species in a given region if
45 there is at least one photograph on Flickr from that region which has been tagged with the
46 name of the species (using either its common name or scientific name). This method is then
47 compared with a standard machine learning based text classification approach, in which all
48 Flickr tags are used, and in which a species may be predicted to occur in a region even if
49 no photographs in that region have been tagged with its name. For the text classifier, we
50 follow the method from [6]. In particular, we show that the best results are obtained by a
51 meta-classifier, which combines the prediction of the text classifier with information about
52 the occurrence of the species name in or near the given region. These results clearly show
53 that better distribution models can be found by taking explicit account of the occurrence of
54 the species name as a tag, in combination with exploiting all other tags.

55 **2** Related Work

56 An overview of the potential for exploiting social media in conservation and biodiversity was
57 provided by Di Mini et al [3], who conducted a study of the use of social media platforms for
58 posting observations of nature. The most commonly used platforms were, in order of level
59 of sharing of nature related content: Facebook, Instagram, Twitter, Youtube, Flickr and
60 LinkedIn. The potential of Flickr for mapping wildlife observations was illustrated by Barve
61 [1] who mapped geotagged postings that included the scientific or common names for the
62 Monarch Butterfly and the Snowy Owl, although that study did not conduct any systematic
63 evaluation of the quality of the retrieved data. Daume [2] performed a manual evaluation of
64 a sample of Twitter postings that named three invasive species (using associated photos for
65 validation). They identified factors correlated with valid observations, such as the presence
66 of a linked photo and tags that describe the environment (e.g. ‘leaves’ and ‘tree’). The
67 present work exploits such associated tags in predicting species distribution. An approach
68 to validating individual observations in Flickr was described by ElQadi et al [5] who used
69 Google’s reverse image-search service to find photos similar to those in Flickr postings. The
70 tags of the Google photos were then compared with those in Flickr in an attempt to filter
71 out non-wildlife images. In our work we learn an association between all Flickr tags and the
72 presence of particular species at a location.

73 The methods presented here build on the work of [6] which exploited weighted values
74 of all tags to train an SVM (support vector machine) classifier to predict the presence of
75 various environmental phenomena including species. In looking at species distribution no
76 distinction was made in [6] between whether the species name was present or not and the
77 focus was on the additional value that Flickr tags provide relative to scientific data such as
78 climate and landcover.

79 **3** Methodology

80 The objective of this paper is to find a method that can use Flickr tags for predicting the
81 occurrence of wildlife species. To this end, we split the target spatial area into grid cells
82 $C = \{c_1, \dots, c_{x_m}\}$ and associate each cell with all the georeferenced Flickr tags that occur

83 within the cell. Following [6], we use Positive Pointwise Mutual Information (PPMI) to
 84 weight how strongly tag t is associated with cell c . In particular, PPMI compares the actual
 85 number of occurrences with the expected number of occurrences (given how many tags
 86 occur overall in c and how common the tag t is). Let $f(t, c)$ be the number of times tag
 87 t (from the set of all tags T) occurs in the cell c . Then the weight $PPMI(t, c)$ is given by
 88 $\max\left(0, \log\left(\frac{P(t, c)}{P(c)P(t)}\right)\right)$ where:

$$89 \quad P(t, c) = \frac{f(t, c)}{N} \quad P(t) = \frac{\sum_{c' \in C} f(t, c')}{N} \quad P(c) = \frac{\sum_{t' \in T} f(t', c)}{N} \quad N = \sum_{t' \in T} \sum_{c' \in C} f(t', c')$$

91 Each cell c is now represented as a sparse vector V_p , encoding the PPMI weight of all the
 92 tags in c . We assume that a training set $K \subset C$ is available which contains cells with known
 93 ground truth species observations and a testing set $U \subset C \setminus K$ containing cells whose species
 94 presence our method will try to estimate.

95 Our method of estimating the presence of a particular species s in cell c involves learning
 96 two classifiers *SVM1* and *SVM2*. The aim of the first classifier *SVM1* is to make initial
 97 predictions for the cells in the testing set U using the feature vector representation V_p . To
 98 give a higher confidence to tags that correspond to the name of the species, we combined the
 99 output of *SVM1* (i.e. classifier confidence score value) with information about the presence
 100 or absence of the *Common Name* or the *Scientific Name* of that species in the cell c or
 101 the neighboring cells. In particular, the cell c is now represented as a feature vector V_m
 102 which contains three features: the confidence value predicted by *SVM1*, the presence of the
 103 species actual name in c as a binary feature (being 1 if the c contains the actual name and
 104 0 otherwise), and the percentage of neighbours that contain the species name (again as a
 105 common or scientific name) as tag. The second classifier *SVM2* is learned using the feature
 106 vector V_m to give the final estimation.

107 **4 Experimental Evaluation**

108 **4.1 Data Acquisition**

109 In this work we use two datasets: the ground truth species distribution from the National
 110 Biodiversity Network Atlas (NBN Atlas)³ and the geocoded social media postings from the
 111 photo sharing website Flickr⁴. The NBN is a collaborative project committed to making
 112 biodiversity information available via the NBN Atlas. This dataset covers the UK and Ireland.
 113 We used the Flickr API to collect approximately 12 million georeferenced Flickr photographs
 114 within the UK and Ireland in September 2015. However, our analysis in this paper will focus
 115 only on the tags associated with these photographs. The NBN Atlas dataset contains a total
 116 of 302 birds with at least 1000 observations, of which 200 have a name that occurs in at least
 117 100 Flickr photographs. Among these, we have considered a random sample of 50 birds for
 118 our experiments. Note that even species with a large number of occurrences may possibly
 119 only occur in a few cells.

120 **4.2 Experimental Settings and Baselines**

121 In the experiments, we consider a binary classification problem for each of the selected birds.
 122 Specifically, the task we consider is to predict in which of the grid cells the bird occurs (i.e. for

³ NBN Atlas occurrence download at <http://nbnatlas.org>. Accessed 19 April 2018.

⁴ <http://www.flickr.com>

<45>:4 Mapping Wildlife Species Distribution

123 which grid cells the NBN Atlas data contains at least one observation). We test our method
124 at three levels of granularity, considering grid cells of size 10, 20 and 30 kilometers. The
125 set of cells C was split into two-thirds for training, one-sixth for testing, and one-sixth for
126 tuning the SVM parameters. It is known that the quality of any supervised model is strongly
127 affected by the way in which the data are divided. Therefore, we split the study area into
128 geographically separated regions, as shown in Figure 1, to test the ability of our method to
129 make predictions about geographic regions for which no observation records are given. This
130 makes the task more challenging than choosing the cells randomly, due to possible differences
131 between the training and testing regions. Finally, for formal evaluation we compared the
132 results of three different methods: “Species Names” which predicts that the species occurs
133 if its common or scientific name appears in at least one Flickr photo in the test cell, “All
134 Flickr Tags” ($SVM1$) which uses the PPMI-based feature vector modelling all Flickr tags
135 to train an SVM classifier using the cells in the training set and predict labels for the cells
136 in the testing cells, and finally “Meta features”($SVM2$) which is our proposed method, as
137 described in Section 3.



■ **Figure 1** Training, Tuning, and Testing regions.

138 4.3 Results and Discussion

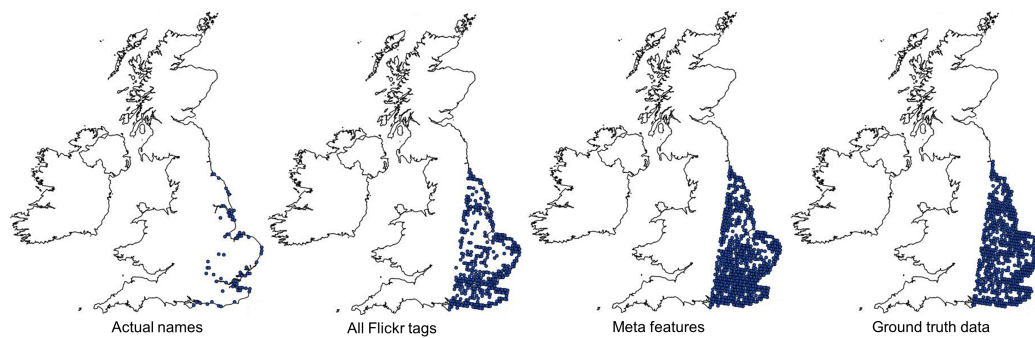
139 The results of predicting species distribution are reported in Table 1 in terms of the average
140 accuracy, average precision, average recall, average F1 score, and average Area Under the
141 ROC Curve (AUC) over the 50 birds. The results clearly show that “All Flickr Tags”
142 significantly outperforms “Species Names”. However, the proposed meta-classifier leads to
143 the best results overall, especially in terms of F1 score.

144 While the “All Flickr Tags” approach works well overall, we found a few cases where
145 using only the species names led to better performance. Perhaps unsurprisingly, this is
146 mostly the case when the number of NBN records (i.e. True labels) in the training region
147 is low, as there may not be enough training data to effectively learn an SVM classifier in
148 such cases. To illustrate such issues, Table 2 shows the F1 scores of 5 individual species.
149 As can be seen, for common species such as Mallard, Dunlin, and Green Sandpiper, the
150 “All Flickr Tags” method performs rather well. In contrast, for some less common species
151 (or species which only occur in particular geographic contexts), such as Atlantic Puffin and
152 Nightingale, we found better results when using the “Species name” method. Interestingly,
153 our proposed meta classifier, which takes account of both the species presence data and the
154 all tags classification for nearby regions, outperforms both of the other methods for almost
155 all the considered species.

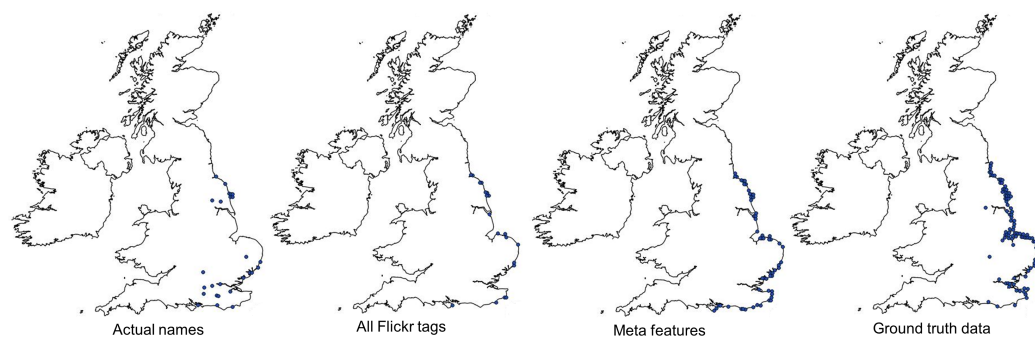
156 Figures 2 and 3 visually illustrate the performance of our method. Note that these species
 157 (like most of the considered birds) occur in fewer than 50% of the cells, which is intuitively
 158 why the “All Flickr Tags” method is more cautious in predicting occurrence (i.e. in absence
 159 of any reason to predict occurrence, it is safer for a classifier to predict non-occurrence).

■ **Table 1** Results for predicting the distribution of 50 species across the testing area.

Dataset	Cell Size	Accuracy	Precision	Recall	F1 Score	AUC
Species Names	10 km	0.520	0.876	0.109	0.183	0.550
All Flickr Tags	10 km	0.779	0.787	0.500	0.560	0.801
Meta features	10 km	0.825	0.820	0.603	0.637	0.850
Species Names	20 km	0.501	0.943	0.241	0.355	0.613
All Flickr Tags	20 km	0.784	0.852	0.639	0.705	0.893
Meta features	20 km	0.870	0.907	0.811	0.832	0.917
Species Names	30 km	0.567	0.970	0.384	0.515	0.684
All Flickr Tags	30 km	0.831	0.868	0.758	0.795	0.943
Meta features	30 km	0.919	0.943	0.896	0.905	0.952



■ **Figure 2** Prediction of the Dunlin distribution across the testing area with 10km grid cells.



■ **Figure 3** Prediction of the Atlantic Puffin distribution across the testing area with 10km grid cells.

160 5 Conclusions and Future Work

161 In this paper we have presented a method for mapping the location of wildlife species
 162 occurrence using the evidence of tags from the photo sharing web site Flickr. We have shown

<45>:6 Mapping Wildlife Species Distribution

■ **Table 2** F1 scores for predicting the distribution of individual species using different methods.

	No.NBN records	No.Flickr photos	Cell size	Species Names	All Flickr Tags	Meta features
Mallard (<i>Anas platyrhynchos</i>)	1718823	11831	10 km	0.640	0.978	0.985
			20 km	0.899	0.974	0.986
			30 km	0.955	0.988	0.992
Dunlin (<i>Calidris alpina</i>)	278872	796	10 km	0.196	0.630	0.744
			20 km	0.346	0.920	0.969
			30 km	0.553	0.980	0.996
Green Sandpiper (<i>Tringa ochropus</i>)	103295	187	10 km	0.077	0.610	0.806
			20 km	0.195	0.849	0.955
			30 km	0.367	0.906	0.980
(Common) Nightingale (<i>Luscinia megarhynchos</i>)	24437	383	10 km	0.128	0.0	0.401
			20 km	0.326	0.0	0.705
			30 km	0.512	0.0	0.835
(Atlantic) Puffin (<i>Fratercula arctica</i>)	11551	2512	10 km	0.152	0.136	0.367
			20 km	0.173	0.359	0.518
			30 km	0.264	0.476	0.630

163 that while a method based simply on the presence or absence of the species name provides
164 good precision, much better overall accuracy, with similar precision, can be achieved with a
165 machine learning classifier that combines the presence-absence data with predictors based on
166 all the textual tags of the photos.

167 One line of future work is to investigate the use of a text classifier to estimate confidence
168 in observations of wildlife species in individual social media postings. This could be of
169 particular value when considering postings that mention a species name but in a context
170 that might be unrelated to its occurrence in nature.

171 — References —

- 172 **1** Vijay Barve. Discovering and developing primary biodiversity data from social networking
173 sites: A novel approach. *Ecological Informatics*, 24:194–199, 2014.
- 174 **2** Stefan Daume. Mining twitter to monitor invasive alien species? An analytical framework
175 and sample information topologies. *Ecological Informatics*, 31:70–82, 2016.
- 176 **3** Enrico Di Minin, Henrikki Tenkanen, and Tuuli Toivonen. Prospects and challenges for
177 social media data in conservation science. *Frontiers in Environmental Science*, 3:63, 2015.
- 178 **4** Janis L. Dickinson, Benjamin Zuckerberg, and David N. Bonter. Citizen science as an
179 ecological research tool: Challenges and benefits. *Annual Review of Ecology, Evolution,
180 and Systematics*, 41:149 – 172, 2010.
- 181 **5** Moataz Medhat ElQadi, Alan Dorin, Adrian Dyer, Martin Burd, Zoe Bukovac, and Mani
182 Shrestha. Mapping species distributions with social media geo-tagged images: Case studies
183 of bees and flowering plants in australia. *Ecological Informatics*, 39:23–31, 2017.
- 184 **6** Shelan S. Jeawak, Christopher B. Jones, and Steven Schockaert. Using flickr for charac-
185 terizing the environment: An exploratory analysis. In *13th International Conference on
186 Spatial Information Theory, COSIT 2017*, pages 21:1–21:13, 2017.
- 187 **7** Philip Lei, Gustavo Marfia, Giovanni Pau, and Rita Tse. Can we monitor the natural
188 environment analyzing online social network posts? a literature review. *Online Social
189 Networks and Media*, 5:51–60, 2018.
- 190 **8** Anthony Stefanidis, Andrew Crooks, and Jacek Radzikowski. Harvesting ambient geospa-
191 tial information from social media feeds. *GeoJournal*, 78(2):319–338, 2013.