

Supporting information for “A scalable machine learning approach for measuring violent and peaceful forms of political protest participation with social media data”

1 Materials and Methods

1.1 Identifying Protest Time and Location with Associated Press (AP) Image Metadata

The adept Bayes classifier built as part of this study was trained on human-coded Tweets labeled with each of the four types of social actions: Singular Peace, Collective Peace, Singular Force and Collective Force. Because the database of 600 million geo-coded Tweets which we had available to us was simply too large for human coders to sift through for the purposes of identifying each of the social actions discussed above, it was necessary for us to narrow down sets of Tweets that would have the highest probability of corresponding to each of the four social actions. While those studying social actions in the past have used searches through media accounts of protest activity there are a number of issues related to using this method which our method of utilizing AP image metadata solve.

First, as discussed above written media accounts of protest activity tend to be biased towards very large scale protests and social actions that often involve violence. Since newspaper databases are mostly in English as well, coverage of protest activity will be biased towards protests occurring in the English-speaking world. AP photographers, on the other hand, cover a much wider variety of protest activity both within English speaking countries and outside of English speaking countries.

This is because many AP photographers work on a freelance basis and will take photographs of events that may or may not be subsequently written about by larger corporate news sources such as CNN, Fox, NBC, CBS and others [1]. Thus, AP images and image metadata related to protest activity will include a much wider breadth and depth of protest and social action activity than will traditional news sources. This, in turn, provides us with a much wider scope for our training data and ultimately a much more useful and accurate classifier.

Second, while newspaper accounts require exhaustive searches through texts to determine the time and place that a social action or protest occurred, Associated Press image metadata contains information about exactly when and where a photo of a particular protest or social action was snapped. Furthermore, all photographs were visually verified as including some type of social action.

As mentioned above we used AP image metadata to identify clusters of Tweets that were more likely to be related to the each of the social actions defined in the paper. To do this we first queried the full AP Image Database which we had subscription access to through the U.C. Berkeley Library by searching for every photo containing the tag “protest” between April 1st, 2014 and April 30th, 2015, the date range that we had geocoded Tweets available for. This search yielded a total of 5,156 photos which were subsequently downloaded by ourselves and 5 U.C. Berkeley undergraduate assistants and compiled in an image database. While we are able to release metadata extracted from these images, Associated Press copyright restrictions prevent us from releasing the database containing the images themselves.

For each image, metadata which included a caption, photographer, date and time

that the photo was taken, country, city and state (where applicable) was available in the metadata.

The Python code which was used to extract and format the image metadata is shown below. Please note that all files in the code below refer to machine readable *.json* versions of the *.csv* data files which we included as part of this submission for purposes of exposition.

```

1 import json, re, os, iptcinfo
2 from iptcinfo import c_datasets, c_datasets_r, IPTCInfo
3
4 root = "~/Dropbox/URAP_APIImages/"
5 outfile = root+"data/events.json"
6
7 people = [
8     "student_1", "student_2", "student_3", \
9     "student_4", "student_5", "student_6", "student_7"]
10
11 events = {}
12 missing = {}
13 for person in people:
14     imagenames = os.listdir(root+person+"/images/")
15     imagenames = [x for x in imagenames if 'jpg' in x]
16     f = open(root+person+"/images/img_info.json", 'w')
17     for image in imagenames:
18         try:
19             iptc = dict(IPTCInfo(image).data)
20             data = {}
21             for numkey in iptc:
22                 data[c_datasets[numkey]] = iptc[numkey]
23             f.writelines(json.dumps(data) + '\n')
24         except:
25             missing.setdefault(person, [])
26             missing[person].append(image)
27     f.close()
28
29     lineNumber = 0
30     try:
31         f = open(root+person+"/images/img_info.json", "r")
32         lineNumber = 0
33         for line in f:
34             lineNumber += 1
35             try:
36                 data = json.loads(line)
37                 city = data["city"]
38                 caption = data["caption/abstract"]
39                 date = data["date created"]
40                 events.setdefault(city, {})
41                 events[city].setdefault("boxes", [])
42                 events[city].setdefault("captions", {})
43                 events[city]["captions"].setdefault(date, [])
44                 if caption not in events[city]["captions"][date]:
45                     events[city]["captions"][date].append(caption)
46             except:
47                 pass
48
49         f.close()
50     except:
51         print person+": no info file!"
52 os.chdir("~/projects/protest-project/")
53 f = open("events.json", "w")
54 f.writelines(json.dumps(events))
55 f.close()
56 f = open("missing_images.json", "w")
57 f.writelines(json.dumps(missing))

```

Fig A Sample of tweets and social action codes in JSON format used for the social action coding task.

```
58 f.close()
```

Code A Python code for AP Image Metadata Extraction

After Code A was run, the collected image metadata held within **events.json** was edited by hand to associate the “city” field to appropriate geographic units within the geo-coded Tweets database, which is annotated with codes for the GADM global administrative boundaries. See <http://www.gadm.org/version2> for more information. The contents of the resulting GADM-associated database of image metadata are dumped into CSV format under the file **APprotestevents-metadata.csv** for convenience.

1.2 Extracting samples of geo-coded tweets Using AP metadata

Using the GADM-associated AP image metadata database **events.json** discussed above, the geo-coded Twitter database was queried and tweets were extracted using geographic and time parameter ranges which corresponded to event data listed in **APprotestevents-metadata.csv**. While alignment of the extracted data to the geo-coded Twitter database was straightforward for temporal information (times were simply converted from GMT to EST), association for geographic data was somewhat more complicated. Generally, the AP images were tagged with city-level spatial information, while the geo-coded Twitter database was tagged for US county-level equivalents. This resulted in each city being coded by hand for overlapping counties, which resulted in different levels of refinement for different cities. For example, while San Francisco was mapped to its name-identical county, New York City was identified with the five boroughs, and each of Oakland and Berkeley were associated with their superset, Alameda county. Ultimately, the effects of this imperfect mapping were most significantly seen at the time of coding the filtered Twitter date for social actions. In particular, for the small cities that are subsets of large, populous counties, like Ferguson MO., the filtering resulted in samples that were more “watered down” in terms of social action (since more unrelated tweets were included from the surrounding areas).

After extraction, tweets were reformatted into a series of JSON formatted files which were then distributed to each of the 5 U.C. Berkeley undergraduate students who were subsequently instructed about how to code each of the tweets that they were assigned according to the four social action categories that we described in the paper. Following this, the student-coded tweets were comprehensively audited for consistency and then merged to form the final coded data set.

A sample of the completed Tweet database as viewed by the coders is presented in Fig A. The Python code which was used to query the Tweet database and format relevant tweets for human coding is below.

```
1 import datetime, os
2
3 f = open("events.json", "w")
4 f.writelines(json.dumps(events))
5 f.close()
6
7 for city in events:
8     dts = []
9     runs = []
10    for d in events[city]:
11        if d != "boxes" and d != "runs":
```

```

12         for t in events[city][d]:
13             ttime = d[0:4] + "-" + d[4:6] + "-" + d[6:8] + "-" + t[0:2] + "-" + t[2:4]
14             dts.append(ttime)
15         dts = sorted(dts)
16         pt = datetime.datetime.strptime(dts[0], '%Y-%m-%d-%H-%M')
17         i = 0
18         runs.append(["NA", "NA"])
19         runs[i][0] = (pt - datetime.timedelta(hours=5.25)).strftime('%Y-%m-%d
20             -%H-%M')
21         for dt in dts:
22             dt = datetime.datetime.strptime(dt, '%Y-%m-%d-%H-%M')
23             if dt - pt > datetime.timedelta(hours=0.5):
24                 runs[i][1] = (pt - datetime.timedelta(hours=4.75)).strftime('%
25                     Y-%m-%d-%H-%M')
26                 i += 1
27                 runs.append(["NA", "NA"])
28                 runs[i][0] = (dt - datetime.timedelta(hours=5.25)).strftime('%
29                     Y-%m-%d-%H-%M')
30                 pt = dt
31             if runs[i][1] == "NA":
32                 runs[i][1] = (pt - datetime.timedelta(hours=4.75)).strftime('%Y-%m
33                     -%d-%H-%M')
34             events[city]["runs"] = runs
35
36         if len(events[city]["boxes"]):
37             boxes = events[city]["boxes"]
38             outfile = "~/projects/protest-project/data/tweets/"+city+".txt"
39             f = open(outfile, "w")
40             f.close()
41             for run in runs:
42                 start = run[0]
43                 end = run[1]
44                 cd = "cd ~/tools/perl/pullTweets; "
45                 cmd = "perl bin/pullTweets.pl -tblucm start="+start+" end="+
46                     end+" boxes="+", ".join(boxes)+" languages=[en] >> \""+
47                     outfile+"\" "
48                 os.system(cd+cmd)

```

Code B Python code using AP image metadata to query geo-coded Tweet database.

The Python code which was used to prepare the tweets for human coding is shown below.

```

1 import re, json
2 import random as ra
3
4 with open("~/projects/protest-project/events.json", "r") as f:
5     events = json.loads(f.read())
6
7 data = {}
8 raw = {}
9 for city in events:
10     infile = "~/projects/protest-project/data/tweets/"+city+".txt"
11     try:
12         f = open(infile, "r")
13         i = 0
14         for line in f:
15             line = line.strip()
16             if re.match("^(.*?)\t(.*?)\t(.*?)\t(.*?)\t(.*?)\t(.*?)$", line):
17                 i += 1
18                 time, box, coords, lang, user, tweet = re.split("\t", line)
19                 ID = city + "-" + str(i).zfill(5)
20                 ID = re.sub("-", "_", ID)
21                 raw[ID] = {
22                     "time": time,

```

```

23         "box": box,
24         "coords": coords,
25         "user": user,
26         "tweet": tweet
27     }
28     data[ID] = {
29         "tweet": tweet,
30         "codes": {
31             "singular-peace": 0,
32             "collective-peace": 0,
33             "collective-force": 0,
34             "singular-force": 0
35         }
36     }
37     f.close()
38 except:
39     print "no data on "+city+"!"
40 f = open("~/projects/protest-project/data/surveys/allTweets-pm15m.json", "w")
41 f.writelines(json.dumps(row))
42 f.close()
43
44 ra.seed(0)
45 tweeKeys = list(data.keys())
46 ra.shuffle(tweeKeys)
47 i = 0
48 j = 0
49 batch = {}
50 for ID in tweeKeys:
51     i += 1
52     batch[ID] = data[ID]
53     if not i % 3000:
54         j += 1
55         f = open("~/projects/protest-project/data/surveys/batches/batch-"+
56                 str(j).zfill(3)+".json", "w")
57         f.writelines(json.dumps(batch))
58         f.close()
59         f = open("~/projects/protest-project/data/surveys/batches/batch-"+
60                 str(j).zfill(3)+".tsv", "w")
61         for ID in batch:
62             line = ID+"\t[0,0,0,0]\t"+json.dumps(batch[ID]["tweet"])+"\n"
63             f.write(line)
64         f.close()
65
66         batch = {}
67         i = 0
68 if len(batch.keys()):
69     j += 1
70     f = open("~/projects/protest-project/data/surveys/batches/batch-"+
71             str(j).zfill(3)+".json", "w")
72     f.writelines(json.dumps(batch))
73     f.close()
74     f = open("~/projects/protest-project/data/surveys/batches/batch-"+
75             str(j).zfill(3)+".tsv", "w")
76     for ID in batch:
77         line = ID+"\t[0,0,0,0]\t"+json.dumps(batch[ID]["tweet"])+"\n"
78         f.writelines(line)
79     f.close()

```

Code C Python code used to prepare tweets for human coding.

1.3 Building the labeled data database for adept Bayes classifier

After coding of the relevant tweets, tweets were compiled into a labeled tweet database containing 22,626 tweets which was used to train the adept Bayes classifier and assess classifier performance. The labeled tweet data used to train and test the classifier is included as part of the supplementary materials in this submission under the filename **Coded-Tweets-training.csv**.

1.4 Building and training the adept Bayes classifier

In addition to the substantive contribution of this paper, we provide a methodological innovation which improves upon a well known machine learning classifier known as naïve Bayes. While the naïve Bayes classifier is known to provide superior performance results and scalability in the context of text classification, a major drawback of the algorithm is the conditional independence assumption which assumes that words within documents are independent of each other. Although the naïve Bayes classifier performs very well for text analysis tasks despite the conditional independence assumption, many have theorized that relaxing this assumption can potentially improve performance even further [2].

Here, we relax the conditional independence assumption by incorporating a NLP phrase chunking method [3] that allows us to better meet the independence assumptions of the naïve Bayes classifier, resulting in an enhanced naïve Bayes classifier which we refer to as “adept” Bayes.

Code for implementing the adept Bayes classifier is shown below.

```

1 import json, re, sys, os
2 import numpy as np
3 import math as ma
4 import random as ra
5 from partitioner.methods import MWE
6 pa = MWE()
7
8 def countPhrases(message):
9     counts = {"*": 0.}
10    for MWE in pa.partition(message):
11        phrase = "".join([block[1] for block in MWE])
12        if re.search("[a-zA-Z]", phrase):
13            counts.setdefault(phrase, 0.)
14            counts[phrase] += 1.0
15            counts["*"] += 1.0
16    return counts
17
18 def loadPhrases(record):
19     codes = record["codes"]
20     counts = countPhrases(record["tweet"])
21     for code in codes:
22         codeCounts = counts
23         if codes[code]:
24             codeCounts["_"+code] = 1.
25             yield code, codeCounts
26         else:
27             codeCounts["_not."+code] = 1.
28             yield "not."+code, codeCounts
29
30 def aggregateByPhrase(d1, d2):
31     for key in d2:
32         d1.setdefault(key, 0)
33         d1[key] += d2[key]
34     return d1
35

```

```

36 def findN(record):
37     ID, counts = record
38     yield "_ALLCODES", counts
39     counts["_N"] = float(len(counts.keys()) - 1)
40     yield ID, counts
41
42 def computeEntropies(record):
43     code, entropies = record
44     n = float(entropies["_N"])
45     entropies["*"] = sum([entropies[phrase] for phrase in entropies\
46                          if phrase != "_N" and phrase != "*" and phrase
47                          != "_" + code])
48
49     m = float(entropies["*"])
50     del entropies["_N"]
51     del entropies["*"]
52     alpha = n/N
53
54     for phrase in entropies:
55         if not re.match("^[^_]", phrase):
56             entropy = -np.log2((entropies[phrase] + alpha) / (alpha * N + m))
57             entropies[phrase] = entropy
58
59     entropies["_DEFAULT"] = -np.log2(alpha / (alpha * N + m))
60     return code, entropies
61
62 def processTweet(tweet):
63     result = {
64         "posteriors": {code: -prior_entropy[code] for code in
65                       prior_entropy},
66         "M_words": 0.,
67         "counts": {},
68         "N": 0.,
69         "M": 0.,
70         "tweet": tweet
71     }
72
73     review_M = 0
74     tempCounts = countPhrases(tweet)
75     for phrase in tempCounts:
76         if phrase != "*":
77             result['counts'].setdefault(phrase, 0)
78             result['counts'][phrase] += tempCounts[phrase]
79
80     todel = []
81     for phrase in result['counts']:
82         if phrase != "*":
83             INLIKE = False
84             for code in likelihoodEntropy:
85                 if likelihoodEntropy[code].get(phrase, False):
86                     result['posteriors'][code] -= result['counts'][phrase]
87                     * likelihoodEntropy[code][phrase]
88                     INLIKE = True
89             if INLIKE:
90                 result['M'] += result['counts'][phrase]
91                 result['M_words'] += result['counts'][phrase] * len(re.
92                               split(" ", phrase))
93                 result['N'] += 1
94                 for code in likelihoodEntropy:
95                     if not likelihoodEntropy[code].get(phrase, False):
96                         result['posteriors'][code] -= result['counts']
97                         ][phrase] * likelihoodEntropy[code]["
98                         _DEFAULT"]
99                 result['counts'][phrase] = int(result['counts'][phrase
100 ])
101             else:
102                 todel.append(phrase)

```

```

95     for phrase in todel:
96         del result['counts'][phrase]
97
98     newPosteriors = {}
99     for code in result['posteriors']:
100         code = re.split("\.", code)
101         if len(code) == 2:
102             code = code[1]
103             newPosteriors.setdefault(code, {})
104             newPosteriors[code]["0"] = result['posteriors'][code]
105         else:
106             code = code
107             newPosteriors.setdefault(code, {})
108             newPosteriors[code]["1"] = result['posteriors'][code]
109     result['posteriors'] = newPosteriors
110
111     for code in result['posteriors']:
112         if len(result['posteriors'][code].keys()) == 2:
113             P0 = 1. / (1. + (2 * (result['posteriors'][code]["1"] - result[
114                 'posteriors'][code]["0"])))
115             P1 = 1. - P0
116             result['posteriors'][code] = {"0": P0, "1": P1}
117         else:
118             result['posteriors'][code] = {"0": 1., "1": 0.}
119
120     result["M"] = int(result['M'])
121     result["M_words"] = int(result['M_words'])
122
123     return result

```

Code D Functional components of the adept Bayes classifier. A complete pipeline uses all but the last function, *processTweet()*, for training, and the final function for classification, as will be seen in the subsequent cross-validation code.

Using the coded data from Twitter, the model was trained as a parallel series of binary classifiers. In other words, a separate adept Bayes classifier was trained for each of the four types of social action (e.g., collective force, singular peace, etc.), in addition to the collapsed categories (i.e., collective, singular, force, peace, or any). This means that the processing of the coded data resulted in the training of 9 binary classifiers that can separately assess the presence of each type of social action, e.g., application of the classifier would predict that a tweet either is a representation of collective force, or is not a representation of collective force, while simultaneously predicting if the same tweet is a representation of singular peace, or is not a representation of singular peace, et cetera.

1.5 Assessing classifier performance

After training the adept Bayes classifier, we assesses precision, recall and F_1 statistics using tenfold cross-validation on the trained model.

The code used to produce the trained model and the classifier statistics are shown below.

```

1 codes = [
2     "collective-force",
3     "collective-peace",
4     "singular-force",
5     "singular-peace",
6     "collective",
7     "singular",
8     "peace",
9     "force",
10    "action"
11 ]

```



```

12 thresholds = [(x)/100. for x in range(1,100)]
13
14
15 records = []
16
17 with open("trainingData.json", "r") as f:
18     for line in f:
19         record = json.loads(line)
20         oldcodes = dict(record["codes"])
21         for code in oldcodes:
22             for newcode in re.split("-", code):
23                 if not record["codes"].get(newcode, 0):
24                     record["codes"][newcode] = oldcodes[code]
25                 if not record["codes"].get("action", 0):
26                     record["codes"]["action"] = oldcodes[code]
27             records.append(record)
28
29 ra.seed(0)
30
31 ra.shuffle(records)
32
33 numfolds = 10
34 folds = []
35
36 foldsize = int((len(records)/numfolds) + 0.5)
37 if foldsize*numfolds < len(records):
38     foldsize += 1
39 allresults = []
40 for foldnum in range(numfolds):
41     trainingrecords = []
42     testingrecords = []
43
44     for k in range(numfolds):
45         if k == foldnum:
46             testingrecords.extend(list(records[k*foldsize:(k+1)*foldsize]))
47         else:
48             trainingrecords.extend(list(records[k*foldsize:(k+1)*foldsize]))
49
50 ##### training
51 #####
52 allcounts = {}
53 for code in codes:
54     allcounts[code] = {}
55     allcounts["not."+code] = {}
56
57 total_tweets = 0.
58 for record in trainingrecords:
59     total_tweets += 1.
60     for k, v in loadPhrases(record):
61         allcounts[k] = aggregateByPhrase(allcounts[k], v)
62
63 for code in codes:
64     allcounts["_ALLCODES"] = aggregateByPhrase(allcounts[code],
65         allcounts["not."+code])
66     break
67 for ky in allcounts.keys():
68     allcounts[ky]["_N"] = float(len(allcounts[ky].keys()) - 1)
69
70 N = float(len(allcounts["_ALLCODES"].keys())) - 2.
71
72 del allcounts["_ALLCODES"]
73
74 entropies = []
75 for ky in allcounts.keys():

```

```

75     entropies.append(computeEntropies((ky, allcounts[ky])))
76
77     likelihoodEntropy = {}
78     prior_entropy = {}
79
80     for x in entropies:
81         code = x[0]
82         entdict = x[1]
83         priorNum = entdict["_"+code]
84         prior_entropy[code] = priorNum
85         del entdict["_"+code]
86         likelihoodEntropy[code] = entdict
87
88     for code in prior_entropy:
89         prior_entropy[code] = -np.log2(prior_entropy[code]/total_tweets)
90     ##### end training
91     #####
92
93     ##### testing
94     #####
95     totals = {}
96     for code in codes:
97         totals[code] = {}
98         for thresh in thresholds:
99             totals[code][str(thresh)] = {
100                 "TP": 0.,
101                 "FP": 0.,
102                 "FN": 0.,
103                 "TN": 0.
104             }
105     for record in testingrecords:
106         result = processTweet(record["tweet"])
107         for code in codes:
108             truth = record["codes"][code]
109             for thresh in thresholds:
110                 st = str(thresh)
111                 if result["posteriors"][code]["1"] >= thresh:
112                     prediction = 1
113                 else:
114                     prediction = 0
115
116                 if truth:
117                     if prediction:
118                         totals[code][st]["TP"] += 1.
119                     else:
120                         totals[code][st]["FN"] += 1.
121                 else:
122                     if prediction:
123                         totals[code][st]["FP"] += 1.
124                     else:
125                         totals[code][st]["TN"] += 1.
126
127     for code in codes:
128         for thresh in thresholds:
129             st = str(thresh)
130             try:
131                 totals[code][st]["TPR"] = totals[code][st]["TP"]/(totals[
132                     code][st]["TP"] + totals[code][st]["FN"])
133             except:
134                 totals[code][st]["TPR"] = 0.
135             try:
136                 totals[code][st]["PPV"] = totals[code][st]["TP"]/(totals[
137                     code][st]["TP"] + totals[code][st]["FP"])
138             except:
139                 totals[code][st]["PPV"] = 0.
140             try:

```

```

139         totals[code][st]["F1"] = (2.*totals[code][st]["TPR"]*
140                                     totals[code][st]["PPV"])/\
141                                     (totals[code][st]["TPR"] + totals
142                                     [code][st]["PPV"])
141     except:
142         totals[code][st]["F1"] = 0
143     try:
144         totals[code][st]["ACC"] = (totals[code][st]["TP"]+totals[
145                                     code][st]["TN"])/\
146                                     (totals[code][st]["TP"]+totals[
147                                     code][st]["TN"]+\
148                                     totals[code][st]["FP"]+totals[
149                                     code][st]["FN"])
147     except:
148         totals[code][st]["ACC"] = 0
149
150     print "the results for fold: "+str(foldnum)
151     for code in codes:
152         st = max([(totals[code][str(thresh)]["F1"], str(thresh)) for
153                   thresh in thresholds], key=lambda x: x[0])[1]
154         totals[code][str(thresh)] = st
155         print code, st, totals[code][st]["PPV"], totals[code][st]["TPR"],
156               totals[code][st]["F1"], totals[code][st]["ACC"]
157     print
158     allresults.append(totals)
159     #####
160     ##### report the F1-optimal posterior threshold for each social action
161     type
162     bestThresh = {}
163     for code in codes:
164         bestThresh[code] = [0, 0]
165         for thresh in thresholds:
166             F1 = np.mean([totals[code][str(thresh)]["F1"] for totals in
167                           allresults])
168             if F1 > bestThresh[code][1]:
169                 bestThresh[code] = [thresh, F1]
170         thresh = bestThresh[code][0]
171         print code, thresh, np.mean([totals[code][str(thresh)]["PPV"] for
172                                     totals in allresults]),\
173               np.mean([totals[code][str(thresh)]["TPR"] for
174                       totals in allresults]),\
175               np.mean([totals[code][str(thresh)]["F1"] for
176                       totals in allresults])

```

Code E Python code used to produce tenfold cross-validation statistics for the adept Bayes classifier.

In addition to this, we sought to assess out-of-domain performance of the classifier for each type of action. To accomplish this, we trained the adept Bayes classifier on all coded data, except for those that were drawn from Hong Kong, which were used for testing. Under this setup, the classifier had no in-domain knowledge of the Hong Kong Democracy protests that were tested upon.

The Python code for the execution of the out-of-domain test is shown below.

```

1 codes = [
2     "singular-peace",
3     "collective-peace",
4     "collective-force",
5     "singular-force",
6     "collective",
7     "singular",
8     "peace",
9     "force",
10    "action"

```

```

11 ]
12
13 thresholds = [(x)/100. for x in range(1,100)]
14
15 trainingrecords = []
16 testingrecords = []
17
18 with open("trainingData.json", "r") as f:
19     for line in f:
20         record = json.loads(line)
21         oldcodes = dict(record["codes"])
22         for code in oldcodes:
23             for newcode in re.split("-", code):
24                 if not record["codes"].get(newcode, 0):
25                     record["codes"][newcode] = oldcodes[code]
26                 if not record["codes"].get("action", 0):
27                     record["codes"]["action"] = oldcodes[code]
28             if re.match("Hong-Kong", record["ID"]):
29                 testingrecords.append(record)
30             else:
31                 trainingrecords.append(record)
32
33 ##### training
34 #####
35 allcounts = {}
36 for code in codes:
37     allcounts[code] = {}
38     allcounts["not."+code] = {}
39
40 total_tweets = 0.
41 for record in trainingrecords:
42     total_tweets += 1.
43     for k, v in loadPhrases(record):
44         allcounts[k] = aggregateByPhrase(allcounts[k], v)
45
46 for code in codes:
47     allcounts["_ALLCODES"] = aggregateByPhrase(allcounts[code], allcounts[
48         "not."+code])
49     break
50 for ky in allcounts.keys():
51     allcounts[ky]["_N"] = float(len(allcounts[ky].keys()) - 1)
52 N = float(len(allcounts["_ALLCODES"].keys())) - 2.
53
54 del allcounts["_ALLCODES"]
55
56 entropies = []
57 for ky in allcounts.keys():
58     entropies.append(computeEntropies((ky, allcounts[ky])))
59
60 likelihoodEntropy = {}
61 prior_entropy = {}
62
63 for x in entropies:
64     code = x[0]
65     entdict = x[1]
66     priorNum = entdict["_"+code]
67     prior_entropy[code] = priorNum
68     del entdict["_"+code]
69     likelihoodEntropy[code] = entdict
70
71 for code in prior_entropy:
72     prior_entropy[code] = -np.log2(prior_entropy[code]/total_tweets)
73 ##### end training
74 #####
75

```

```

76 ##### testing
77 #####
78 totals = {}
79 for code in codes:
80     totals[code] = {}
81     for thresh in thresholds:
82         totals[code][str(thresh)] = {
83             "TP": 0.,
84             "FP": 0.,
85             "FN": 0.,
86             "TN": 0.
87         }
88 for record in testingrecords:
89     result = processTweet(record["tweet"])
90     for code in codes:
91         truth = record["codes"][code]
92         for thresh in thresholds:
93             st = str(thresh)
94             if result["posteriors"][code][1] >= thresh:
95                 prediction = 1
96             else:
97                 prediction = 0
98
99             if truth:
100                 if prediction:
101                     totals[code][st]["TP"] += 1.
102                 else:
103                     totals[code][st]["FN"] += 1.
104             else:
105                 if prediction:
106                     totals[code][st]["FP"] += 1.
107                 else:
108                     totals[code][st]["TN"] += 1.
109 for code in codes:
110     for thresh in thresholds:
111         st = str(thresh)
112         try:
113             totals[code][st]["TPR"] = totals[code][st]["TP"] / (totals[code][st]["TP"] + totals[code][st]["FN"])
114         except:
115             totals[code][st]["TPR"] = 0.
116         try:
117             totals[code][st]["PPV"] = totals[code][st]["TP"] / (totals[code][st]["TP"] + totals[code][st]["FP"])
118         except:
119             totals[code][st]["PPV"] = 0.
120         try:
121             totals[code][st]["F1"] = (2.*totals[code][st]["TPR"]*totals[code][st]["PPV"]) / (totals[code][st]["TPR"] + totals[code][st]["PPV"])
122         except:
123             totals[code][st]["F1"] = 0.
124         try:
125             totals[code][st]["ACC"] = (totals[code][st]["TP"] + totals[code][st]["TN"]) / (totals[code][st]["TP"] + totals[code][st]["TN"] + totals[code][st]["FP"] + totals[code][st]["FN"])
126         except:
127             totals[code][st]["ACC"] = 0.
128
129 ## picking up here!
130 st = max([(totals[code][str(thresh)]["F1"], str(thresh)) for thresh in thresholds], key=lambda x: x[0])[1]
131 totals[code]["thresh"] = st
132 print code, st, totals[code][st]["PPV"], totals[code][st]["TPR"], totals[code][st]["F1"], totals[code][st]["ACC"]

```

133 #####

Code F Python code used to produce statistics from out-of-sample validation results.

1.6 Exploring social actions during New York City's climate change protest.

As part of this submission we also include a movie file which demonstrates how our trained adept Bayes classifier can track social actions as they unfold in real-time on the ground. This was accomplished by using the classifier to classify tweets during a multiple day climate change protest in New York City which began on September 21st, 2014. The movie file included as part of this submission is a Quicktime *.mov* file entitled **NYC-ClimateChangeProtests-Med.mov**.

The code and classified tweets used to produce this movie file and are included as part of this submission. The classified tweets used to create this movie file and S2 Fig 2 in the paper is **NY-CliMarch-classified.csv**.

The Python code which was used to create Figures 1, 2 and 3 in the paper is included below.

```
1 from pylab import *
2 import json, re, os
3 import matplotlib.pyplot as py
4 import matplotlib
5 from datetime import datetime
6 import smopy
7 from haversine import haversine
8 import numpy as np
9 from matplotlib.patches import Wedge
10 from geopy import geocoders
11
12 with open("~/projects/protest-project/data/tweets/timeseries/allTimes.json", "r") as f:
13     allTimes = json.loads(f.read())
14
15 with open("~/projects/protest-project/data/tweets/timeseries/timeseries/allActions/New York.timeseries.json", "r") as f:
16     timeseries = json.loads(f.read())
17
18 thresholds = {
19     "collective-force": 0.08,
20     "collective-peace": 0.78,
21     "singular-force": 0.92,
22     "singular-peace": 0.85,
23     "collective": 0.79,
24     "singular": 0.87,
25     "peace": 0.88,
26     "force": 0.71,
27     "action": 0.93
28 }
29
30 gn = geocoders.GeoNames(username = "username")
31
32 outdir = "~/projects/protest-project/data/tweets/timeseries/frames/processed/allActions/"
33 county = ["New York", "NY"]
34 handle = county[0]
35 figdir = "~/Documents/talks/APSA/images/"
36
37 colors = {
38     "singular-peace": "black",
39     "collective-peace": "blue",
40     "collective-force": "red",
```

```

41     "singular-force": "green",
42 }
43
44 codes = [
45     "singular-peace",
46     "collective-peace",
47     "collective-force",
48     "singular-force",
49     "collective",
50     "singular",
51     "peace",
52     "force",
53     "action"
54 ]
55 likelihoodEntropy = {}
56 for amb in codes:
57     f = open("./binary/likelihood-entropy/Alameda_11-24_
58         _firstNight_recoded_stitched_allActions_"+amb+".json", "r")
59     for line in f:
60         d = json.loads(line)
61         for code in d:
62             likelihoodEntropy[code] = d[code]
63
64 ## for NY climate change protests
65 start = "2014-09-21-06-41"
66 end = "2014-09-27-09-50"
67 fignum = 1
68 for hour in timeseries["hours"]:
69     fig = py.figure(fignum, figsize=(24, 16))
70     fig.suptitle(", ".join(county)+" : "+re.sub(":00", "", str(datetime.
71         strptime(hour, "%Y-%m-%d-%H").strftime("%m/%d/%Y %I:%M %p"))),
72         fontsize=25, fontweight='bold')
73     first_ax = fig.add_axes([0.0, 0.75, 1.0, 0.25])
74     fignum += 1
75     ylimits = [0, 0]
76     xlims = [datetime.strptime(start, "%Y-%m-%d-%H-%M"), datetime.strptime(
77         end, "%Y-%m-%d-%H-%M")]
78     xpt = datetime.strptime(hour, "%Y-%m-%d-%H")
79     if xpt < xlims[0] or xpt > xlims[1]:
80         continue
81     for action in colors.keys():
82         x = []
83         y = []
84         for i in range(0, len(timeseries["hours"])):
85             ypt = float(timeseries["posteriors"][action][i])
86             xpt = datetime.strptime(timeseries["hours"][i], "%Y-%m-%d-%H")
87             if xpt >= xlims[0] and xpt <= xlims[1]:
88                 if ypt > ylimits[1]:
89                     ylimits[1] = ypt
90                 if timeseries["hours"][i] == hour:
91                     ix = i
92                     y.append(ypt)
93                     x.append(xpt)
94
95     py.plot(x, y, color = colors[action], lw = 2, label = re.sub("-", " ",
96         action))
97     py.plot(x[ix], y[ix], color = colors[action], ms = 25, marker = ".")
98
99     py.ylabel('Posterior presence of action', fontsize=25)
100     py.xticks(fontsize = 20, rotation = 30)
101     py.yticks(fontsize = 20)
102     py.ylim(ylimits)
103     py.xlim(xlims)
104
105     py.legend(loc='upper right', fontsize = 25)

```

```

102     ### now for the map #####
103     second_ax = fig.add_axes([0.0, 0.0, 0.6, 0.675])
104
105     py.xticks(fontsize = 0)
106     py.yticks(fontsize = 0)
107
108     hoi = hour
109     doi = re.sub("-\\d\\d$", "", hoi)
110     code = "collective-force"
111
112     infile = outdir+handle+"_"+re.sub("-\\d\\d$", "", hoi)+".json"
113
114     f = open(infile, "r")
115     for line in f:
116         record = json.loads(line)
117         if record["hour"] == hoi:
118             frame = record
119             break
120
121     dims = (8,12)
122
123     meanx, meany = gn.geocode(", ".join(county))[1]
124
125     lon = [meanx - 0.03, meanx + 0.14]
126     lat = [meanx - 0.03, meanx + 0.08]
127
128     smomap = smopy.Map((lat[0], lon[0], lat[1], lon[1]), z=14)
129     smomap.show_mpl(figsize=dims, ax = second_ax)
130
131     pts = []
132     pxpts = []
133
134     for cluster in frame["clusters"]:
135         xs, ys = smomap.to_pixels(np.asarray(cluster["center"][1]), np.
136             asarray(cluster["center"][0]))
137         if len(pts) < 2:
138             pts.append([cluster["center"][1], cluster["center"][0]])
139             pxpts.append([xs, ys])
140         else:
141             break
142     try:
143         hdist = haversine(pts[0], pts[1])
144         pxdist = (((pxpts[0][0] - pxpts[1][0]) ** 2) + ((pxpts[0][1] -
145             pxpts[1][1]) ** 2)) ** 0.5
146         ppm = pxdist/(hdist*1000.)
147     except:
148         ppm = 1.
149     fig = py.gcf()
150
151     for cluster in frame["clusters"]:
152         xs, ys = smomap.to_pixels(np.asarray(cluster["center"][1]), np.
153             asarray(cluster["center"][0]))
154         circ = py.Circle((xs,ys), ppm*(cluster["radius"]+25.), color='blue',
155             alpha = 0.5)
156         fig.gca().add_artist(circ)
157
158     ## determine wedge angles!
159     ang = 360.*(float(cluster["positives"][code])/float(cluster["size"]
160         ))
161     if ang:
162         sect = Wedge((xs,ys), ppm*(cluster["radius"]+25.), 0., ang,
163             color="red", alpha=1.)
164         fig.gca().add_artist(sect)
165
166     mostIntense = 0

```



```

162 intenseCluster = {}
163 points = []
164 for cluster in frame["clusters"]:
165     intensity = float(cluster["positives"][code])
166     for point in cluster["points"]:
167         if point["posteriors"][code]['1'] > thresholds[code]:
168             points.append(point)
169     if intensity > mostIntense:
170         intenseCluster = cluster
171     mostIntense = intensity
172 cluster = intenseCluster
173
174 ### now for the shift #####
175 third_ax = fig.add_axes([0.635, 0.0, 0.365, 0.675])
176
177 counts = {}
178 for point in points:
179     for phrase in point["counts"]:
180         counts.setdefault(phrase, 0.)
181         counts[phrase] += point["counts"][phrase]
182
183 phrases = counts.keys()
184 for phrase in phrases:
185     if likelihoodEntropy[code].get(phrase, False) and \
186        likelihoodEntropy["not."+code].get(phrase, False):
187         counts[phrase] *= (likelihoodEntropy["not."+code][phrase] \
188                           - likelihoodEntropy[code][phrase])
189     elif likelihoodEntropy[code].get(phrase, False) and not \
190          likelihoodEntropy["not."+code].get(phrase, False):
191         counts[phrase] *= (likelihoodEntropy["not."+code]["DEFAULT"] \
192                           - likelihoodEntropy[code][phrase])
193     elif not likelihoodEntropy[code].get(phrase, False) and \
194          likelihoodEntropy["not."+code].get(phrase, False):
195         counts[phrase] *= (likelihoodEntropy["not."+code][phrase] \
196                           - likelihoodEntropy[code]["DEFAULT"])
197     else:
198         del counts[phrase]
199 bardata = sorted(zip(counts.keys(), counts.values()), key=lambda x: abs
200                  (x[1]), reverse=False)
201
202 vals = [dat[1] for dat in bardata[len(bardata)-40:len(bardata)]]
203 labs = [dat[0] for dat in bardata[len(bardata)-40:len(bardata)]]
204 pos = arange(len(vals))+.5 # the bar centers on the y axis
205
206 barh(pos, vals, align='center', color = "gray")
207 for i in range(len(vals)):
208     if vals[i] >= 0:
209         text(vals[i]+.01*max(map(abs, vals)), pos[i]-.4, labs[i],
210              fontsize=17)
211     else:
212         text(vals[i]-.01*max(map(abs, vals)), pos[i]-.4, labs[i],
213              fontsize=17, ha = "right")
214 ylabel('Rank by impact', fontsize=15)
215 xlabel('Relative impact on classification', fontsize=15)
216 xticks(fontsize = 15)
217 yticks(fontsize = 15)
218 ylim([-0.5, len(vals)+0.5])
219 try:
220     xlim([min(vals)-.15*max(map(abs, vals)), max(vals)+.35*max(map(abs
221     , vals))])
222 except:
223     continue
224
225 labels = map(str, sorted(range(1,41), reverse=True))
226 py.yticks(pos, labels, rotation='horizontal')

```

```

223
224     plot_margin = 0.05
225
226     x0, x1, y0, y1 = py.axis()
227     py.axis((x0 - plot_margin*100,
228             x1 + plot_margin*100,
229             y0 - plot_margin*2,
230             y1 + plot_margin*2))
231
232     py.savefig(figdir+handle+'_3-panel_'+hoi+'.png', format = "png",
233               bbox_inches='tight', pad_inches=0.25)
234     py.close(fig)

```

Code G Version of the Python code used to create NYC climate march figures.

1.7 Copyright and Compliance

In collecting this data and producing this document, we have complied with all of the data collection terms and conditions stipulated by the *Associated Press (AP) Images Database* hosted by the University of California, Berkeley Library System and with all of the data collection terms and conditions stipulated as stipulated by *Twitter Inc.*.

All map images provided in the document were produced by the *Mapbox* package in *Python*: <https://github.com/mapbox/mapbox-sdk-py>. Tiles produced by the *Mapbox* package can be used under the Creative Commons 3.0 license (CC BY 3.0). More about this can be found here: <http://maps.stamen.com>.

Any questions or concerns regarding these matters should be directed to:

L. Jason Anastasopoulos

Email: ljanastas@uga.edu

References

1. Ricchiardi S. Covering the World: As US news organizations have backed away from foreign news coverage, the Associated Press' international report has become increasingly vital. *American Journalism Review*. 2007;29(6):32–40.
2. Rennie JD, Shih L, Teevan J, Karger DR, et al. Tackling the poor assumptions of naive bayes text classifiers. In: *ICML*. vol. 3. Washington DC); 2003. p. 616–623.
3. Williams JR, Bagrow JP, Reagan AJ, Alajajian SE, Danforth CM, Dodds PS. Selection models of language production support informed text partitioning: an intuitive and practical, bag-of phrases framework for text analysis. *arXiv preprint arXiv:160107969*. 2016;.