

A Data Mining Framework for Analysing Geospatial-Temporal Data

Raunak Sarbajna
Department of Computer Science
Lamar University
Beaumont, Texas, USA
Email: rsarbajna@lamar.edu

Sujing Wang
Department of Computer Science
Lamar University
Beaumont, Texas, USA
Email: sujing.wang@lamar.edu

Abstract—In this paper, we introduce a new data mining framework for analyzing spatio-temporal data. It introduces polygon based change analysis techniques and automatic storytelling methodology for spatio-temporal data. Change analysis and automatic storytelling are essential techniques in understanding larger patterns and trends in multifaceted, timeseries geographic data. We evaluate the effectiveness of our framework through case studies involving Twitter emotion data and North American Drought data. The experimental results show that our framework can discover interesting patterns and useful information from spatial-temporal data.

Index Terms—change analysis, storytelling, sentiment analysis, polygon, spatio-temporal data.

I. INTRODUCTION

Analyzing change in spatial data is critical for many applications including developing early warning systems that monitor environmental conditions, detecting political unrest and crime monitoring.

Change analysis models are essential in understanding larger patterns and trends in multifaceted, time-series geographic data. The purpose of this study is to detect spatiotemporal changes in land use within sequential (time-series) maps. Changes in land use can be categorized by the complex interaction of structural and behavioural factors associated with technological capacity, demand, and social relations that affect both environmental capacity and the demand, along with the nature of the environment of interest.

Modern technology digitizes wide sources of information constantly, with hour after hour of data being stored and most of it being unprocessed raw information. This is especially true with remote sensing networks and resultant geo-spatial imagery, coming from satellites, drone captures and flyover shots. Processed, georeferenced datasets from sources as diverse as gravity measurement a la GRACE [23] to ocean circulation mapping [24].

The goal of this research project is to detect and analyze how the patterns of features change over time and space in spatiotemporal land use datasets. All polygons under consideration are closed spatial georeferenced sets, rather than raster imagery.

Our approach provides a change monitoring framework which creates a change graph that captures the changes in spatial land uses clusters and a change summarization framework

that creates specific change summaries based on the change graph based on the change story types.

Our research contributions are summarized as follows:

- 1) A novel change analysis framework that works on vectored polygonal datasets
- 2) New change predicates that are data agnostic and can work on a large spectrum of data
- 3) A new measure of interestingness to aid in generating automated storytelling based the change analysis results

The rest of the paper is structured as follows. Section 2 reviews previous literature on the subject and discusses related work. Section 3 introduces our data mining framework. Section 4 lays out the methodology in detail. Section 5 evaluates the framework with case studies on drought datasets pre and post Harvey, pre and post California wildfires and twitter emotion polygons. Section 6 provides a conclusion and discusses potential future expansions to the framework.

II. RELATED WORK

A survey of the classical change detection algorithms can be found in the Lu et al. [3] paper and tells us that the integrated GIS and remote sensing approaches yield the best results. However, they are very sensitive to registration accuracies between images. Thus, images must be properly orthorectified and georeferenced, especially because the changes in the emotion polygons are so subtle. This assumes the emotions are to be treated as just another feature in the map, like any other category.

Since our data is primarily in an urban environment, with all the grid like rigidity that entails, it is a good idea to look at change detection algorithms optimized for urban environments. One of the hardest aspects to measure is to distinguish between change and no-change, as well as different kinds of change. Comparing image differencing, image regression, tasseled-cap transformation and chi square transformation, Ridd and Liu [3] find image differencing to be the most consistent, with a sustained overall accuracy of >80%.

It is useful to have a programming-oriented study comparing several of the change detection algorithms using MATLAB, rather than pure application-oriented comparison, in order to have a benchmark. Minu and Shetty [5] analyzed image differencing, image ratioing, change vector analysis, tasseled

cap transformation and principal component analysis for efficiency and effectiveness. Although their area of study was not urban but a variety of land use/ land cover, change vector analysis gave the best overall accuracy. We also studied two novel methods that are recent developments and are showing promising results: Neighborhood Correlation Image and Comprehensive Change Detection Method, both of which are optimized for remote sensing imagery but can be adapted to vectorized maps without loss of generality.

The change detection model using Neighborhood Correlation Image (NCI) logic works because of the obvious fact that the same geographic area (e.g., a 3x3 pixel window) on two dates of imagery will tend to be highly correlated if little change has occurred, and uncorrelated when change occurs [1]. Computing the piecewise correlation between two data sets demonstrates that NCIs contain change information and that NCIs may be powerful tools for change detection.

A high-performance remote sensing method called Comprehensive Change Detection Method (CCDM) integrates spectral-based change detection algorithms and a novel change model called Zone, which extracts change information from two Landsat image pairs [2]. This can be easily modified to work on the Twitter-based emotional grading maps. This method is simple, easy to operate, widely applicable, and capable of capturing anthropogenic changes like our area of interest.

Storytelling techniques are effective summarization method to succinctly organize extensive information. Traditional storytelling has been mostly successful on news articles, blogs, as well as structured databases. However, traditional storytelling techniques tend to perform poorly on social media content, such as Twitter, where text lacks proper form and function [11]. Moreover, the ability to support dynamic storylines as they evolve is critical to modeling fast moving social media streams such as Twitter. Dos Santos et al. [21] introduced a set of methods to automatically derive stories over linked entities in tweets. They model a story as a graph of entities propagating through spatial regions in a temporal sequence, and controls search space complexity by suggesting regions of exploration. They developed algorithms to conduct storytelling to model tweets over space and time, reasoning over spatio-temporal features, and devise spatio-temporal storylines based on connectivity strength.

Kumar et al. [14] proposed an efficient storytelling implementation that embeds the CARTwheels [15] redescription mining algorithm which utilizes induced classification trees to model redescription in an A* search procedure, using the CARTwheels to supply next move operators on search branches to the A* search procedure. Vocht et al. [15] proposed the implementation of an optimized algorithm controlling the pathfinding process to obtain more homogeneous search domain and retrieve more links between adjacent hops in each path to improve the semantic relatedness of concepts mentioned in a story by increasing the relevance of links between nodes through additional domain delineation and refinement steps. Chen et al. [20] proposed a multimodal

imitation learning via generative adversarial networks (MIL-GAN) method to directly model users' interests as reflected by various data by imitating users' demonstrated storylines. MIL-GAN model is designed to learn the reward patterns given user-provided storylines and then applies the learned policy to unseen data. Santos et al. [21] combined storytelling and Spatio-logical Inference (SLI) to generate rules of interaction among entities and measure how well they forecast a real-world event.

Hossain et al [13] introduced Google Fusion Tables(GFT) that offers collaborative data management in the cloud for data scientists to enable the integration of increasingly complex geospatial data to support storytelling. The paper focused on introduction of overview of map processing in GFT, the architecture overview of GFT, and how to scale to large datasets, massive and complex polygon datasets. GFT provides a useful tool for storytelling through interactive maps.

Kumar et al. [14] formulated storytelling as a generalization of redescription mining. Stories are defined as chains of redescriptions. They proposed an efficient storytelling algorithm as A* search around the outputs of a CARTwheels redescription mining algorithm. The efficiency and scalability of the proposed algorithm were evaluated by three application case studies: word overlaps in large English dictionaries, exploring connections between gene sets in a bioinformatics data set, and relating publications in the PubMed index of abstracts.

Hossain et al. [19] proposed an approach to automatically construct stories between entities in large document collections that can help from directed chains of relationships, with support for co-referencing, evidence marshaling, and imposing syntactic constraints on the story generation process. A new optimization techniques based on concept lattice mining is used to rapidly construct stories on massive datasets.

Chen et al. [20] introduced an approach, multimodal imitation learning via generative adversarial networks (MIL-GAN) for generating storyline on unseen data. It can directly model users' interests as reflected by various data. This approach is used to learn the reward patterns given user-provided storylines and then applies the learned policy to unseen data.

Santos et al. [21] introduced three methods of association analysis, Distance-based Bayesian Inference, Spatial Association Index, and Spatio-logical inference, to capture relatedness among real-world events in high data volumes, and to model similar events that are described disparately under high data variability. It takes as input a set of geotemporally-encoded text streams about violent events called "storylines". This study demonstrated that spatio-temporal storytelling is able to capture important associations among violent events reported in social media and traditional datasets.

III. METHODOLOGY

A. Point Data Sources

Our initial approach to this problem was to store all shapefiles in a postgres database with a GIS addon and perform operations in python. We used psycopg2 and osgeo libraries to import, process and visualize maps. However, this lead to

many problems with interconversions between georeferencing schemes, while converting from WKT geometry to PostGIS geography.

We start with basic point data, which contains latitude/-longitude, along with metadata identifying value of interest, whether that is drought level or emotion value. We insert the contents of the shapefile into a PostGreSQL database using the shp2pgsql toolkit that comes along with the PostGIS extension. The results can be seen in Figure 1.

id	avg	long	lat	id	drought	emotion	numtwts	centroid	batchnm	geodata	geom
integer	numeric	numeric	numeric	integer	numeric	numeric	integer	numeric	integer	character varying (80)	geometry
1	1	-78.85166700000000	42.49361100000001	1	0.000300000000000	0.888500000000000	1	0.030024489606718	1	WGS84	01010000...
2	2	-78.85793050000002	42.45197875999999	2	0.000300000000000	0.114286666666667	30	0.030024489606718	1	WGS84	01010000...
3	3	-78.85859999999999	42.45374270000001	2	0.000300000000000	0.114286666666667	30	0.030024489606718	1	WGS84	01010000...
4	4	-78.87389900000003	42.47902016000001	2	0.000300000000000	0.114286666666667	30	0.030024489606718	1	WGS84	01010000...
5	5	-78.93611948999999	42.45111015000001	2	0.000300000000000	0.114286666666667	30	0.030024489606718	1	WGS84	01010000...

Fig. 1. Twitter shapefile data inserted into Postgres db

As our framework relies on using polygonal data, we cannot use this. So we begin by creating a convex hull of points and inserting it into a new table, while preserving related metadata. Due to engine limitations, we need to be careful and insert only those convex hulls that are polygons specifically, not points or lines.

- 1) Our query for this purpose was:

```
INSERT INTO public.june1poly (avgscor ,
numtwts , geodata , id , batchnm , geom)
```

```
(SELECT d.avgscor , d.numtwts , d.geodata ,
d.id , d.batchnm , ST_ConvexHull
(ST_Collect(d.geom))
FROM public."2014-06-01 " AS d
GROUP BY (d.id , d.avgscor , d.numtwts ,
d.geodata , d.batchnm)
```

```
HAVING ST_GeometryType(ST_ConvexHull
(ST_Collect(d.geom))) = 'ST_Polygon')
```

- 2) Then we insert the centroid of each polygon into the table using the query:

```
UPDATE public.june1poly
SET centroid=ST_Centroid(geom)
```

- 3) We repeat this process for every shapefile needed.

Next we run our change predicates, which include:

- 1) $S - Continuing(c, m) \leftrightarrow Agreement(c, m) \geq 0.8$
- 2) $B - Continuing(c, b) \leftrightarrow Oap(c, b) \geq 0.8$
- 3) $Growing(c, m) \leftrightarrow Containmment(c, m) \geq 0.9$
- 4) $Shrinking(c, m) \leftrightarrow Growing(m, c)$
- 5) $Disappearing(c) \leftrightarrow \exists i(belong - to(c, i))$
- 6) $Novel(c) \leftrightarrow \exists i(belong - to(c, i) \text{ and } (i = 1 \text{ or not } (B - Continuing(c, i - 1)))$
- 7) $Shifting$

Which are defined as:

- $Agreement(c, m) = (area(c \cap m)) / (area(c \cup m))$
- $Overlap(s, f) = area(p \cap (p_1 \cap \dots \cap p_m)) / area(p)$

We will demonstrate three of these predicates here.

- 1) To detect polygons that are increasing in size, we check for similar IDs, intersection and then the rate of overlap. We initially check whether the polygons intersect at all before querying for amount of overlap. This leads to faster processing as it discards the many combinations where the polygons don't touch each other. Our query is structured as:

```
SELECT DISTINCT j2.*
FROM public.june1poly j1 ,
public.june2poly j2
WHERE ST_INTERSECTS(j1.geom , j2.geom)
AND
(ST_AREA(ST_INTERSECTION(j2.geom , j1 ))
/ st_area(j2.geom)) > .85
```

- 2) To detect polygons that are shrinking in size, we check for similar IDs, and lower rates of overlap. This can be modified based on need.

```
SELECT DISTINCT j2.*
FROM public.june1poly j1 ,
public.june2poly j2
WHERE ST_INTERSECTS(j1.geom , j2.geom)
AND
(ST_AREA(ST_INTERSECTION(j2.geom , j1 ))
/ st_area(j2.geom)) < .25
```

- 3) To detect polygons that have shifted we compare their centroids and check if they have moved over 75km. Remember these polygons are created through a convex hull of points, which cannot ensure the centroid will lie within the polygon itself. Which is why we are taking a sufficiently large bounding value for the polygon.

```
SELECT ST_Distance_Spheroid
(j1.centroid , j2.centroid ,
'SPHEROID["WGS 84",6378137,298.25]'),
j1.id FROM public.june1poly j1 ,
public.june2poly j2
WHERE j1.id = j2.id AND
ST_Distance_Spheroid(j1.centroid ,
j2.centroid ,
'SPHEROID["WGS 84",6378137,298.25]') > 75000;
```

A sample of what these generated maps look like can be found in Figure 2

B. Polygonal Map Data Sources

Our next, more successful approach was done through three primary set operations: union, intersection and erase. We calculate area of each individual polygon within each map layer. We then execute a union operation and calculate area. The union layer now contains the original areas of both layers and the areas of the overlapping polygons - we now need to query them properly to prepare for calculating the change percentage and tabulating intersection.

To outline the polygon, we examine several different methods:

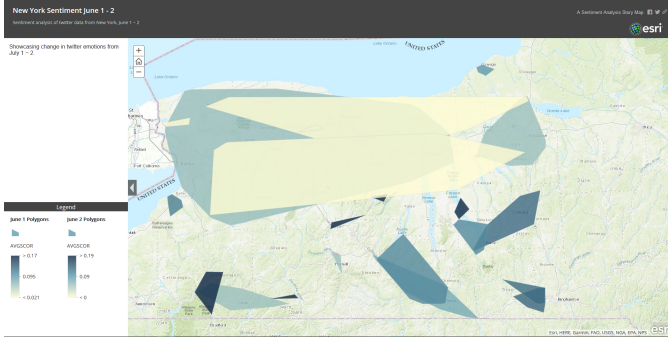


Fig. 2. Polygons generated from Sentiment analysis

- 1) We find features common to either of the layers but not both, essentially performing a symmetrical difference
- 2) We erase the larger of the polygons from the smaller, thus retaining only the growth, and do vice-versa for shrinking
- 3) We perform simple intersection and then invert selection to get changed regions.

Our approach then combined several techniques:

- 1) **Data Pre Processing** This involves curation of datasets with obvious georeferencing errors. This would preferable be done by minimizing the root mean square error. We initially.
- 2) **Parametrization of polygons** Calculate shape and area parameters for each individual polygon with each map layer.s
- 3) **Analysis through Symmetrical Difference** Extract features common to either of the layers
- 4) **Polygon Union Computation** Union sequential layers to contain the original areas of both layers and the areas of overlapping polygons
- 5) **Polygon Erase Operation** Erase larger polygons from smaller (or vice-versa) for detection of growth/shrinking
- 6) **Polygon Intersection/Invert** Selecting and then labelling the changed regions

Our original data source for the drought sets were defined by the USDA as seen in Figure 3.

IV. CHANGE STORYTELLING

In order to able to tell a coherent spatio-temporal data story from the change analysis output, we need to be able to pick out resultant change polygons that have seem to have significant impact on the dataset as a whole. We propose using an interestingness function to look for these.[11]

We choose a set of change polygons SCP. These not only contains those objects but also their associated characteristics. For example, a SCP could contain a set of spatial clusters with polygon, their average drought score, total area, centroid coordinates and other summaries for each spatial cluster. We define the function as:

$$f : SCP \rightarrow [0, \infty)$$

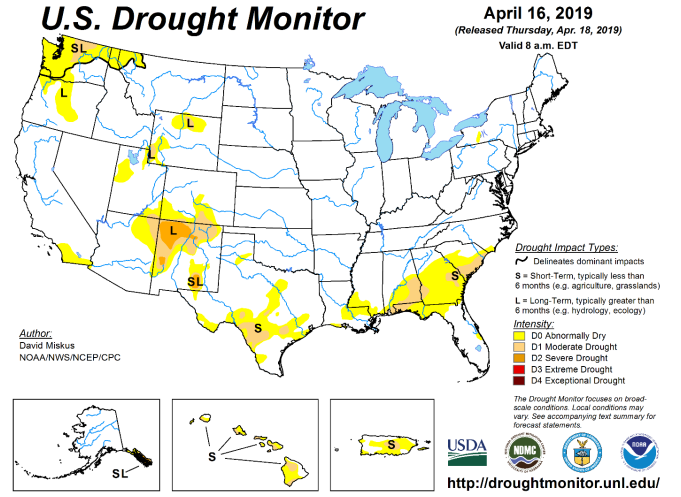


Fig. 3. United States Drought Monitor Data

We define a threshold ω , which ensures that a narrative will only be generated an an object $p \in SCP$ such that $f(p) \geq \omega$. Example parameters for ω include

- $Max \left(\frac{area(Polygon P_i)}{\sum_1^i area(Polygon P_i)} \right)$
- $Max(Percentage\ Change\ in\ Polygon)$
- Largest shift in polygon centroids

The threshold parameters need to be finely tuned so as to not exclude those polygons who fall through exceptions. Once we have a suitable selection of polygons and have chose a threshold value, we can create a summary narrative.

V. CASE STUDY

A. Data Sources

There are two different datasets under consideration here. First is the Spatiotemporal Drought Datasets from the North American Portal: <https://www1.ncdc.noaa.gov/pub/data/nidis/shapefiles/>.

The dataset was constructed using "Geotagged Twitter posts from the United States" [9] and the tool used for creating the emotion clusters was the K2 framework [7]. The raw dataset contains the timestamp, longitude, latitude and the text of each tweet, which are then processed and tokenized.

According to K2, the best package for the emotion score was the Valence Aware Dictionary and Sentiment Reasoner (VADER) system [10], whose analyser parses the tokenized text and checks within a lexicon for words with strong sentiment. The final score is created from the weighted average of all sentimental words and lies within the range [-1, 1], as is required for our Aconagua implementation.

The spatial reference for the Drought shapefiles can be seen in figure 1, and for the twitter data can be seen in figure 2.

The layer definition for the data sets can be seen in tables 1 and 2.

```
GEOGCS["WGS 84",
  DATUM["WGS_1984",
    SPHEROID["WGS 84",6378137,298.257223563,
      AUTHORITY["EPSG","7030"]],
    AUTHORITY["EPSG","6326"]],
  PRIMEM["Greenwich",0,
    AUTHORITY["EPSG","8901"]],
  UNIT["degree",0.0174532925199433,
    AUTHORITY["EPSG","9122"]],
  AUTHORITY["EPSG","4326"]]
```

Fig. 4. Spatial Reference for Drought data

```
GEOGCS["GCS_WGS_1984",
  DATUM["WGS_1984",
    SPHEROID["WGS 84",6378137,298.257223563]],
  PRIMEM["Greenwich",0],
  UNIT["Degree",0.017453292519943295],
  AUTHORITY["EPSG","4326"]]
```

Fig. 5. Spatial Reference for Drought data

B. Emotion Spatial Clusters from Twitter

We are implementing and improving upon the change analysis framework called Aconcagua [6]. The system expects an input of emotion polygons annotated with emotion assessment scores, with +1 representing a very high positive emotion and -1 representing a very high negative emotion. While this method does lead to inconsistencies in locations due to georeferencing inaccuracies, we find that the 8 10m precision [8] works well within city limits.

There are several ways of using the numerous point data we have obtained from this step into an actual polygonal map:

- 1) Creating closed contour lines for contour lines that lie on the boundary of the observation area.
- 2) Creating a convex hull from points with similar scores.

We elaborate on creating convex hulls in the following section.

The original polygons for the drought datasets, before our analysis was started, can be found in Figure 7

We focus on two specific regions to highlight. First, figure 8 shows the regions where areas of drought grew in California following the wildfires in 2017. We noticed patterns

TABLE I
LAYER SPECIFICATION FOR DROUGHT DATA

Name	Type	Width	Precision
long	Real	24	15
lat	Real	24	15
id	Integer	9	0
dnstyTh	Real	24	15
avgScor	Real	24	15
numTwts	Integer	9	0
stdDev	Real	24	15
batchNm	Integer	9	0
geoData	String	80	0

TABLE II
LAYER SPECIFICATION FOR TWITTER DATA

Name	Type	Width	Precision
FIPS_ADMIN	String	4	0
GMI_ADMIN	String	7	0
ADMIN_NAME	String	42	0
FIPS_CNTRY	String	2	0
GMI_CNTRY	String	3	0
CNTRY_NAME	String	40	0
POP_ADMIN	Integer	9	0
TYPE_ENG	String	26	0
TYPE_LOC	String	50	0
SQKM	Real	16	2
SQMI	Real	16	2
COLOR_MAP	String	2	0

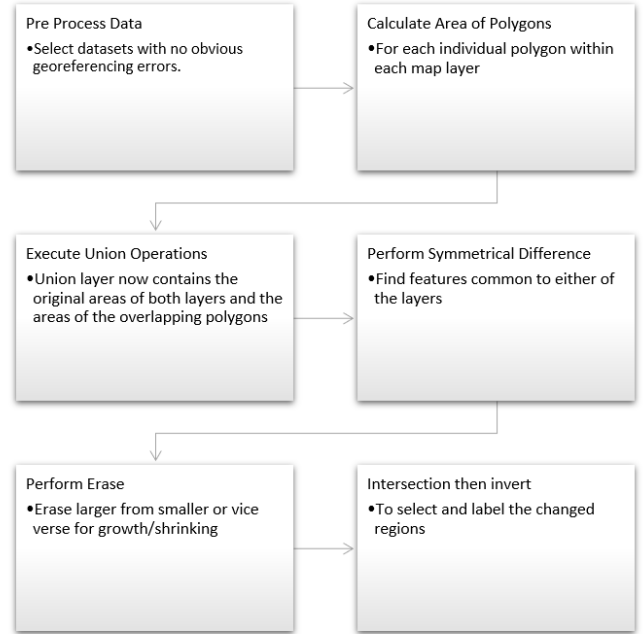


Fig. 6. The Framework Architecture

of increasing drought surrounding the regions that were burnt down, with especially large period upstream from rivers. Our procedure worked well dealing with the large number of small polygons in the region that occur due to isolated wildfires. However, we had difficult dealing with convex polygons that intersected with each other multiple times.

Next, we inspect the region in Texas after Hurricane Harvey. Figure 9 shows the regions where drought affected regions increased following the disaster and Figure 10 shows the regions where they decreased.

We observe that drought affected regions decrease at a high rate around the South Eastern Texas and Louisiana region, which follows common logic. However, there is no clear relation between drought prone regions and river basins close to the coast. We believe that high amount industrial regions create a micro climate that affects the water content entering the soil.

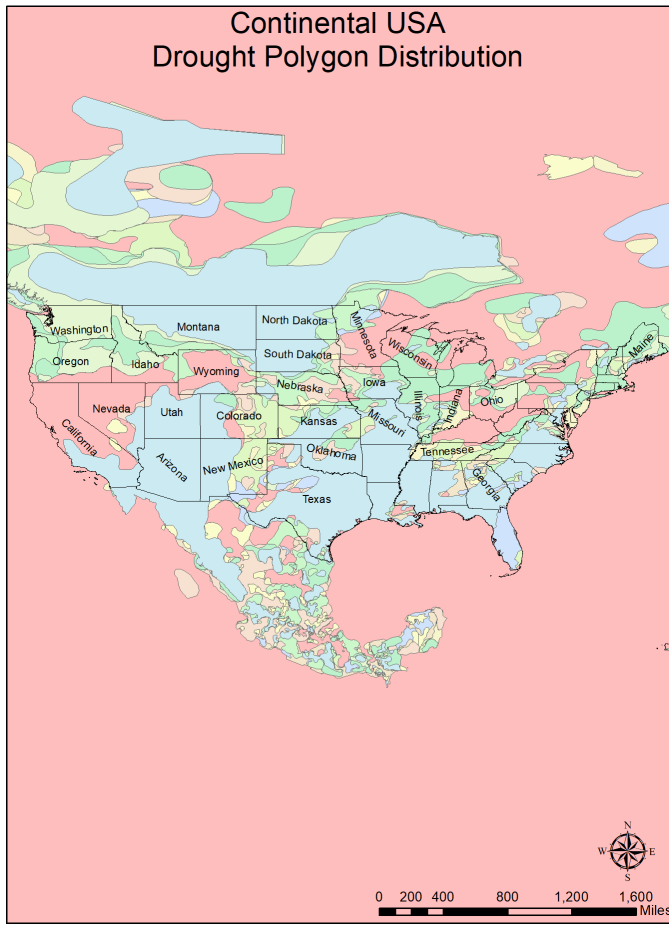


Fig. 7. Continental USA Original Drought polygon Distribution for 2017

According to our results, the drought prone regions increased substantially around the West Texas region both during and after Hurricane Harvey, with a larger increase as winter came around. However, this could be an artefact of the fact that the original dataset had very few polygons in the Texas area, which leads to broader conclusions. Our technique is still severely dependent on the resolution of the input imagery.

The results from the analysing the twitter emotion maps can be seen in Figure 11.

VI. CONCLUSION AND FUTURE WORK

Our experimental studies show that our change detection and analysis framework can successfully detect changes in both polygonal and point-set land use spatiotemporal datasets.

There is a lot work left in dealing with irregularly shaped convex polygons that appear in poorly georeferenced real world data. We believe this is because the sensitivity of the geographic operations in PostGIS or ArcGIS. We are currently working on extending our framework to deal with edge cases like those.

We are also looking towards consolidating and extending some of the change predicates, especially trying to tune them for fine changes.

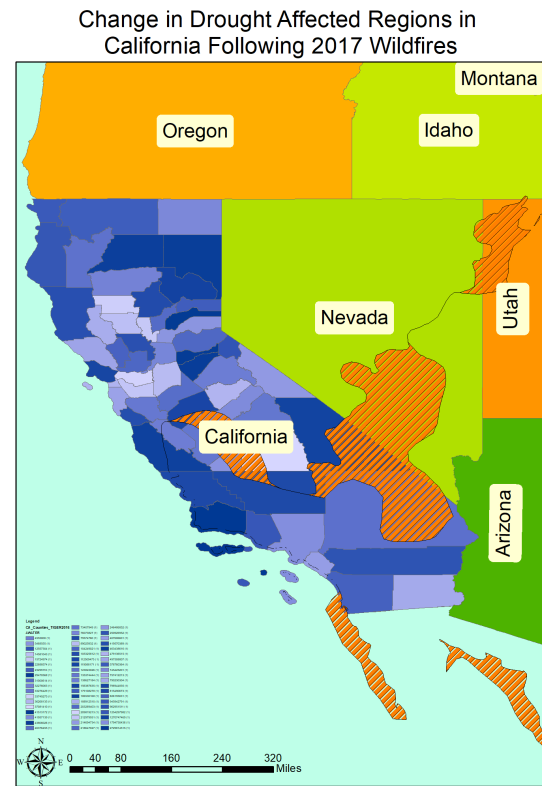


Fig. 8. Change in drought affected regions in California following 2017 wildfires.

We will work on the creating change summaries and stories, as detailed in Section IV, from the data in future.

ACKNOWLEDGMENT

This research is supported by Lamar University Research Enhancement Grant

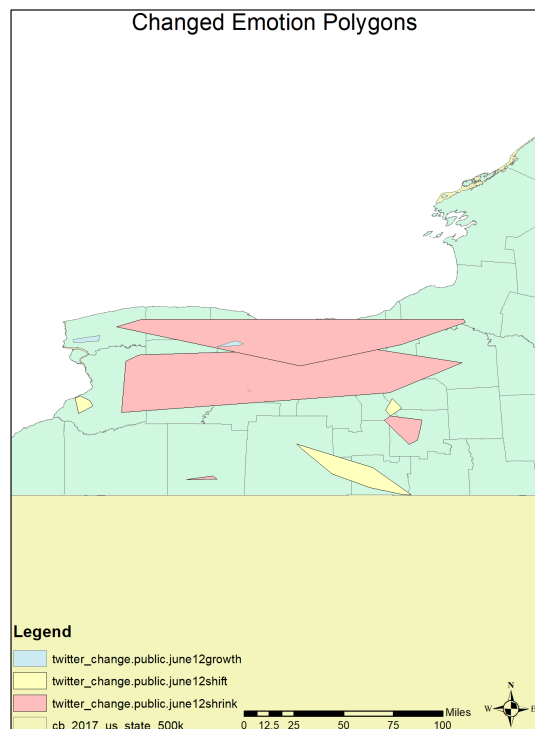


Fig. 11. Change in emotion Polygons in New York State

- Proceedings of the 26th International Joint Conference on Artificial Intelligence, August 19-25, 2017, Melbourne, Australia.
- [17] R. D. Santos, S. Shah, F. Chen, A. Boedihardjo, Chang-Tien Lu, N. Ramakrishnan, Forecasting location-based events with spatio-temporal storytelling, in Proceedings of the 7th ACM SIGSPATIAL International Workshop on Location-Based Social Networks, pp. 13-22, November 04-04, 2014, Dallas/Fort Worth, Texas.
 - [18] Jayant Madhavan, Sreeram Balakrishnan, Kathryn Brisbin, Hector Gonzalez, Nitin Gupta, Alon Halevy, Karen Jacqmin-Adams, Heidi Lam, Anno Langen, Hongrae Lee, Rod McChesney, Rebecca Shapley, Warren Shen, "Big Data Storytelling through Interactive Maps"
 - [19] M. Shahriar Hossain¹, Patrick Butler¹, Arnold P. Boedihardjo², Naren Ramakrishnan, Storytelling in Entity Networks to Support Intelligence Analysts, KDD '12 Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining, Pages 1375-1383, Beijing, China, August 12 - 16, 2012
 - [20] Zhiqian Chen, Xuchao Zhang, Arnold P. Boedihardjo, Jing Dai, Chang-Tien Lu, Multimodal storytelling via generative adversarial imitation learning, Proceedings of the 26th International Joint Conference on Artificial Intelligence, August 19-25, 2017, Melbourne, Australia
 - [21] Raimundo F. Dos Santos, Jr., Arnold Boedihardjo, Sumit Shah, Feng Chen, Chang-Tien Lu, Naren Ramakrishnan, The big data of violent events: algorithms for association analysis using spatio-temporal storytelling, Geoinformatica, v.20 n.4, p.879-921, October 2016
 - [22] Deept Kumar, Naren Ramakrishnan, Richard F. Helm, and Malcolm Potts. 2008. Algorithms for Storytelling. IEEE TKDE20, 6 (2 2008), 736-751. DOI: <http://dx.doi.org/10.1145/1188913.1188915>
 - [23] "GRACE Tellus Data," GRACE Tellus. [Online]. Available: <https://grace.jpl.nasa.gov/data/get-data>. [Accessed: 26-Jul-2019].
 - [24] "SARAL - eoPortal Directory - Satellite Missions." [Online]. Available: <https://directory.eoportal.org/web/eoportal/satellite-missions/s/saral>. [Accessed: 26-Jul-2019].