**Earthquake Prediction and Mitigation Through the Use of Data Based Anomalies**

Jordan Clarke, Daniel Hernandez, Christopher King, Oli Negeri, Daniel Yu-Cua

University of Maryland Global Campus

DATA 610 Decision Management Systems Spring 2022

Dr. Laila Moretto

April 5, 2022

**Introduction**

The data set used in this analysis is the earthquake data set, acquired from The CORGIS Dataset Project at https://think.cs.vt.edu/corgis/csv/earthquakes. This data set has 8,394 rows and 17 columns of data containing earthquake records collected from the United States Geological Survey from July 27, 2016, until August 25, 2016, as shown in Figure 1. The variables in this data set used to describe each earthquake are distance, gap, magnitude, depth, and significance. Additional variables are used for the exact geological coordinates and description of the earthquake's location.

This data set has 4 variables that were determined to be irrelevant and 2 variables that have some erroneous data. The variable id is a unique ID for each earthquake in the data set. This variable was determined to be irrelevant because it does not have a relationship with any of the other variables. The variable hour is the hour that the earthquake occurred. This variable was determined to be irrelevant because it is duplicate information of the variable time.full contains the full time and date the earthquake occurred. The variable minute is the minute that the earthquake occurred. This variable was determined to be irrelevant because it is duplicate information of the variable time.full contains the full time and date the earthquake occurred. The variable second is the second that the earthquake occurred. This variable was determined to be irrelevant because it is duplicate information of the variable time.full contains the full time and date the earthquake occurred.

The variable gap is the largest azimuthal gap in degrees between the azimuthally adjacent stations. The azimuthal gap ranges from 0 to 180. The variable magnitude is a measure of the size of an earthquake at its source. The magnitude is a logarithmic measurement, a magnitude 5 earthquake is 10 times larger than a magnitude 4 earthquake. The magnitude ranges from -1 to

10. The variable significance is a number that describes how significant the earthquake is. The significance is determined by magnitude, maximum MMI, felt reports and estimated impact. The significance ranges from 0 to 1000. The variable depth is the depth of the earthquake in kilometers. The variable distance is the distance the earthquake occurred from the reporting station measured in degrees. The distance ranges from 0.4 to 7.1. The variable location.full is the full name of the location of the earthquake. The variable latitude is the degree of latitude the earthquake occurred. The latitude ranges from -90 to 90 where negative numbers are for southern latitudes. The variable longitude is the degree of longitude the earthquake occurred. The longitude ranges from -180 to 180 where negative numbers are for western longitude. The variable name is the state or country where the earthquake occurred. The variable time.full is the full date and time for when the earthquake occurred. The variable day is the day of the month the earthquake occurred. The variable month is the month when the earthquake occurred. The variable month was changed from a numeric value to text as shown in Figure 2. The variable year is the year when the earthquake occurred.

Split variable location.full to extract the city and state into the variable City/State then split that variable to extract the city into the variable City where the earthquake occurred. Renamed the variable name to State/Country. Split the variable time.full to extract the time the earthquake occurred into the variable time and to extract the day of the week the earthquake occurred into the variable Day of the Week.

Upon initial evaluation of the data set the variables gap and distance contained some erroneous data. The variable gap was filtered to only contain numbers between 0 and 180 as shown in Figure 3. The variable distance was filtered to only contain numbers between 0.4 and 7.1 as shown in Figure 4 1,747 of the 8,394 earthquakes were removed which left 6,647

earthquakes to be analyzed.

**Data Exploration**

Data provided shows a large range of earthquakes that vary in significance and magnitude from disastrous to practically unnoticeable. The earthquakes recorded over this period do not allow for a proper prediction of future earthquakes as the period does not even scale the entire year. However, it does open the floor to other possible industries such as insurance, construction, and government law. For construction, it would be the research and development of methods and practices to construct more durable and stabler buildings in high-risk areas. In the case of insurance, it would be locating, zoning, and determining insurance claims based on various factors such as home durability and potential earthquake risk. Finally for government law, it would be the case of mandating those newer construction practices and limiting insurance discrimination based on the findings and explorations of the dataset.

The first visualization shows that 74.8% of the earthquakes occurred in Alaska and California as shown in Figure 5. If there can be a proper identification of these "high-risk areas" zoning and areas can be identified to notify local home and business owners of the potential risk of these areas. Another idea would be to superimpose the tectonic plate lines onto the first figure to allow for a better understanding of the underlying plates. Next, it's important to understand whether the significance of an earthquake is correlated with its magnitude. While magnitude is used to determine the amplitude of the seismic waves the significance will ultimately be used to determine the impact of the earthquake. Graphing the two side by side as shown in Figure 6 does not indicate a direct relationship. A slight adjustment of the exploration as shown in Figure 7 yields an extremely visually clear example of how they are related to one another. The comparison can also be done when it comes to depth and how it can be related to significance as

shown in Figure 8. Seeing that the two have some loose correlation can reassure us that these two factors will be reliable predictors of the significance of the impact. If scientists can measure the depth and the magnitude, we can then have a rough estimate of what the significance might be.

The visualization shown in Figure 9 was created to understand which state/country had the highest magnitude by average. The visualization only displays the top ten highest magnitude by average. By doing this visualization we have learned that in August, Burma had the highest average of 5.77 and in July, Greece had the highest average of 5.2. When looking at the top ten highest averages in both months, the only State/Country which appears on both diagrams is 'South Georgia and the South Sandwich Islands" with an average of 5.05 in July and 4.98 in August.

The visualization shown in Figure 10 was created to understand which states within the United States with the highest magnitude by average. This visualization was filtered to all states in the USA and displays only the top ten highest averages in July and August. The insights gained from this visualization are that in July, the state of Georgia had the highest average of 4.2 and in August the state of Oklahoma had the highest average. The state that appears in both months is the state 'Puerto Rico' with an average of 2.04 in July and 2.31 in August.

The visualization shown in Figure 11 is represented in terms of magnitude, significance, and depth. It would be useful to check the influence of depth, significance, and magnitude on an earthquake. This visualization is to answer the question of what the data of earthquakes look like monthly. Here the data given of earthquakes are visualized based on the magnitude, significance, and depth monthly. The average value of depth recorded ranges from a minimum of 22.72 in August to a maximum of 28.92 when the month is July. The average most deep earthquake in August. The average values of significance recorded range from a minimum of 56.56 in August

to a maximum of 61.05 when the month is July. August is the most frequently occurring category of the month with a count of 5,467 items 82.2 % of the total. The average values of magnitude recorded range from a minimum of 1516 in August to a maximum of 1601 when the month is July. August is the slightly most occurring category of the month with a count of 1601 items 52 % of the total. In all measurements of the earthquake data recorded, the higher average data is recorded in August.

**Decision Tree**

Two decision trees were generated to evaluate the drivers and generate rules to predict the magnitude as shown in Figure 12 and distance as shown in Figure 13 of earthquakes as measured from the various detection sites around the globe. Magnitude is measured using the Richter scale and provides a scale of reference for seismic activity in terms of damage, lower values are minor earthquakes and damage is minimal if present and larger numbers indicate high levels of damage caused by the earthquake measured Distance is a significant variable to analyze as it was the strongest single driver of the variables not filtered out due to applicability and is a measurement of the approximate distance of the earthquake from the detection site in degrees (e.g., 1° of latitude is approximately 111 km). Variables that were removed from the analysis and decision tree were day, month, and significance. Day and month were removed because the data set is from August 2016, any analysis based on the date would be lacking comparisons from other years. Significance was removed as during the initial review of the data set and variables, significance was identified as being a number generated based on other data within the set, such as magnitude, and estimated MMI (Modified Mercalli Intensity, another form of measuring impact). Due to these factors, significance had a single variable predictor score of 77%, more than double the nearest other drivers.

Reviewing the decision tree and rules aimed at predicting magnitude as shown in Figure 14, we find that distance is the main factor in determining magnitude, followed by depth and gap. According to the tree, one should expect to see the strongest earthquakes at distances greater than .15° (16.7 km) and depths greater than 10 km. According to the predictors, the earthquakes that occur at depths of 20km or higher the gap (a measure of uncertainty in the location of the event measured between detection sites) can be used to predict the magnitude. An example of a prediction using the rules generated by the model shows that an earthquake with a magnitude of 3.82 would occur at a distance greater than .15° and 20 km below sea level or deeper.

Given the relevance to the measured magnitude that we observed with distance, the second decision tree model we ran used distance as the target, with the same excluded variables of day, month, and significance. The predictive model was much simpler, with magnitude acting as the main driver. According to the model, earthquakes with magnitudes below 1.99 had an average distance of .07°; but there were no additional rules. For earthquakes with magnitudes greater than 1.99, the model uses depth and gap respectively. Based on a review of the rules produced by this model as shown in Figure 15, we feel that the magnitude predictive model is our preferred model due to the spread and availability of data to drive the model. Distance proves to be very useful in predicting the measured magnitude of seismic events, however, based on a comparison of the two models the inverse does not appear to be the same. The application of a rules-based approach in an organization based on the types of rules produced by these models is various: profit or loss leaders could be identified, cost drivers can be isolated, marketing budget or direction can be ordered based on transactional areas, etc. Implementation of said rules-based approach would require organizational support, relevant data, investment into analytical tools, and implementation of processes that support and are driven by the rules.

**Dashboard**

This first dashboard tab highlights the most active regions as shown in Figure 16. The

dashboard also provides the average magnitude, depth, and distance. It also shows the maximum

recorded magnitude for the top 10 locations. By focusing on one of the top locations by activity,

or by using a filter within the dashboarding tool, one can drill down to the relevant metrics for a

given region with ease. Another feature is the quick reference of magnitude on the graph using

the color of each bar, as can be seen in the first image California and Alaska combined had

approximately 75% of all the earthquakes during the period from July 27, 2016, until August 25,

2016. California had 1,180 more earthquakes than Alaska did, however, the earthquakes that

occurred in Alaska had a higher average magnitude of 1.49 than California at 1.06. The strongest

earthquake occurred in the Northern Mariana Islands at 7.7 magnitudes while the weakest

average. Organizations use this type of dashboard to include reviewing sales or marketing

statistics (can be targeted by region, date, business unit, product line, etc.), quality control

review/audit, loss leader analysis, and goal setting/tracking.

For the second dashboard tab as shown in Figure 17, we decided to focus on earthquakes

that were reported within the United States region to see how the data compares to the rest of the

regions. This dashboard displays which top and bottom five states reported the highest and

lowest magnitude were felt during the period of this dataset. Alaska had the highest magnitude

with 5.6 and the lowest magnitude belonged to New Jersey. The average magnitude of

earthquakes reported in the United States is 1.22 with Alaska having the highest magnitude

average with a magnitude of 5.6 and North Carolina having the lowest average magnitude of

0.95. The average depth of earthquakes reported in the United States is a depth of 16.23km with

the lowest depth being 3.39km above sea level in Utah. While the largest depth was found in

Virginia at 17.17 km. Another interesting fact what that 18% of all of the earthquakes that occurred in the United States occurred on a Thursday while only 12.4% occurred on a Tuesday.

The third dashboard tab as shown in Figure 18, shows the state and or country where earthquakes occurred during the period from July 27, 2016, until August 25, 2016, for each of the 4 hemispheres. There was a total of 6,114 earthquakes that occurred in the Northwestern Hemisphere. California had the most earthquake occurrences with 3,075 earthquakes while 13 different states or countries had only one occurrence of an earthquake. There was a total of 6,114 earthquakes that occurred in the Northwestern Hemisphere. California had the most earthquake occurrence with 3,075 earthquakes while 13 different states or countries had only one occurrence of an earthquake. The total number of earthquakes in the Northeastern Hemisphere was 250. The Northern Mariana Islands had the most earthquake occurrences with 52 earthquakes while 11 different countries had only one occurrence of an earthquake. There was a total of 133 earthquakes in the Southeastern Hemisphere. Indonesia had the most earthquake occurrences with 43 earthquakes while 6 different countries had only one occurrence of an earthquake. The total number of earthquakes in the Southwestern Hemisphere was 150. Chile had the most earthquake occurrences with 33 earthquakes while 7 different countries had only one occurrence of an earthquake. Each map also displays a heat map of the locations of each earthquake. The Northwestern Hemisphere has most of its earthquakes on the west coast of the United States and along the Alaskan coast. The Northeastern Hemisphere has most of its earthquakes along a line from eastern Russia down to Japan and then down to the Mariana Islands. The Southeastern has most of its earthquakes along the southern border of Indonesia east through Papua New Guinea through the Solomon Islands and to the east on New Caledonia. The Southwestern Hemisphere has most of its earthquakes along the western coast of Peru and Chile with several also on the

east coast of the South Georgia and South Sandwich Islands.

Magnitude and significance both describe the overall power and scale of an earthquake as shown in Figure 19, but it's important to understand the distinction when determining the importance of earthquakes in a particular area. In the dataset, magnitude is described as "a measure of the size of an earthquake at its source. It is a logarithmic measure. At the same distance from the earthquake, the amplitude of the seismic waves from which the magnitude is determined is approximately 10 times as large during a magnitude 5 earthquake as during a magnitude 4 earthquake. The total amount of energy released by the earthquake usually goes up by a larger factor; for many commonly used magnitude types, the total energy of an average earthquake goes up by a factor of approximately 32 for each unit increase in magnitude. Typically ranges from -1 (very tiny) to 10 (incredibly powerful)." This capped level of intensity indicates that it is part of the Richter scale. The Richter scale was introduced in 1935 by Charles F. Richter and Beno Gutenberg. The calculation is done by calculating the logarithm of the amplitude of the largest seismic wave. This scale has been replaced by the Moment Magnitude scale which was introduced in the late '70s and provided an unbound upper bound to the earthquakes (Rafferty, J. P., and Pletcher, Kenneth (2022, March 16)).

While magnitude describes the size of the wave significance describes the actual impact that the earthquake has on the public. The actual equation that describes the calculation of significance is unavailable to us, the description states that significance is a combination of factors such as magnitude, maximum MMI, felt reports, and estimated impact. MMI is short for "the Modified Mercalli Intensity which estimates the shaking intensity from an earthquake at a specific location by considering its effects on people objects, and buildings" (Association of Bay Area Governments Resilience Program. (2013)).

An example of how Magnitude and significance are correlated but not directly causative is the comparison between the major Haiti earthquake of 2010 and the major earthquake of Japan in 2011. It is estimated that roughly 250,000 people died in the Haiti earthquake while only 20,000 were killed in Japan. The magnitude of Haiti's earthquake was 7.0 while Japan's was 9.0 (Miyamoto, H. K., Gilani, A. S. J., &amp; Wong, K. (2011)). Both received an MMI rating of X at their epicenters but had significantly different impacts on the community. An important factor to consider was that Haiti's earthquake epicenter was centered in the country while Japan's was on its coastline with the tsunami doing most of the damage (Hanazato, T. (2011, March)).

The dashboard in the figure is to view and analyze the trend of the earthquake significance, magnitude, and depth monthly based on the data presented. The data covers the month of July and August of 2016. The variables are discussed separately and analyzed the trend is combined at the end. First, one of the influencing parameters on data is the depth of the earthquake. The data shows that the decline combined as far as the depth is concerned from August to July. Earthquake influences in August are in a high position both on average and summation of the depth of the earthquake.

Secondly, to analyze the effect of the magnitude of the earthquake. The trend of magnitude data is much less in July compared to August. A mechanism of a high magnitude earthquake is more trustable and has bold meaning in terms of the severity of the earthquake. The earthquake is more likely to be represented by magnitude to be true to get attention.

Another measure of the relative strength of an earthquake is the significance of which the shaking is noticed. This measure has been particularly useful in estimating the relative severity of shocks that were recorded for the analysis in the data presented. In July, the significance registered was low compared to that of August. The average strength of the earthquakes recorded

in all parameters is higher in August than in July. Finally, the dashboard helps to analyze the strength of the earthquake at a particular time in terms of the significance, depth, and magnitude combined.

**Story**

For the story, the focus was on the earthquake activity in the United States separated into regions Eastern, Central, and Western. The following link will provide access to the video of the story presentation https://youtu.be/BnX9eNo9wV8.

**Conclusion**

As we all know not only the amount of data is important, but also the quality of data. The azimuthal gap is denoted as the largest angle between stations recording an earthquake. A gap that is 180 degrees and over is considerably less accurate than a narrow gap. By pinpointing hotspots with high gap numbers, we can systematically place seismic stations in places that need better quality data. The housing industry is another major market that can greatly benefit from earthquake analytics. Some homeowners' and renters' insurance do not cover earthquake damage, so it is very important to keep in mind what kind of coverage they may need. For instance, in a high-risk area such as San Francisco "assuming coverage for a median home price of $1.3 million, rates could be close to $6,000 a year" (Bishop, L.). Analyzing data such as this can help homeowners better understand the risks of buying a home in a high-risk state such as San Diego or Los Angeles. It also helps insurance companies better predict and zone for potential risks propagated into the future.

# References

Association of Bay Area Governments. (2013). *Making Sense of the Modified Mercalli Intensity*

*Scale (MMI) – A Measure of Shaking*.

https://abag.ca.gov/sites/default/files/making_sense_of_the_modified_mercalli_intensity_

scale.pdf

Bishop, L. (2022, January 27). *How Much Does Earthquake Insurance Cost in California?*

ValuePenguin. https://www.valuepenguin.com/california-earthquake-insurance-cost

Bytenskaya, Y. (2021a, March 14). *Cognos Analytics - Story Properties* [Video]. YouTube.

https://www.youtube.com/watch?v=UyR154W69-

Y&list=PL2r2WGYKOnJWvwjzgazpro6MzdmvLZqBj&index=22

Bytenskaya, Y. (2021b, March 14). *Cognos Analytics - Export dashboard as a story* [Video].

YouTube.

https://www.youtube.com/watch?v=an93WQsmyH8&list=PL2r2WGYKOnJWvwjzgazpr

o6MzdmvLZqBj&index=21

Bytenskaya, Y. (2022, March 3). *Cognos Analytics Session 2 Spring 2022* [Video]. YouTube.

https://www.youtube.com/watch?v=bgimvsEL-z8

Communications and Publishing. (2014, February 24). *The 1964 Great Alaska Earthquake and*

*Tsunami*. United States Geological Survey. https://www.usgs.gov/news/state-news-

release/1964-great-alaska-earthquake-and-tsunami

Cox, A., & Hart, R. B. (1991). *Plate Tectonics*. Wiley.

Hanazato, T. (2011, March 27). *Tohoku Pacific Earthquake on 11 March 2011*. International

Council on Monuments and Sites.

https://www.icomos.org/risk/2011/ICOMOS_Japan_%20201103_earthquake_reports_20
110331.pdf

Miyamoto, H. K., Gilani, A. S. J., & Wong, K. (2011). Massive Damage Assessment Program
and Repair and Reconstruction Strategy in the Aftermath of the 2010 Haiti Earthquake.
*Earthquake Spectra*, *27*(1_suppl1), 219–237. https://doi.org/10.1193/1.3631293

Rafferty, J. P., & Pletcher, K. (2022, March 16). *Japan earthquake and tsunami of 2011*.
Encyclopedia Britannica. https://www.britannica.com/event/Japan-earthquake-and-
tsunami-of-2011

Tiira, T., Uski, M., Kortström, J., Kaisko, O., & Korja, A. (2015). Local seismic network for
monitoring of a potential nuclear power plant area. *Journal of Seismology*, *20*(2), 397–
417. https://doi.org/10.1007/s10950-015-9534-8

# Appendix A

## Figure 1

*Earthquake Data Sample*

| gap ▼ | magnitude | significance | depth | distance ▼ | latitude | longitude | City | State/Country | Date | year | month | day | Day of the week | Time |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ | ↑↓ |
| 122 | 1.43 | 31 | 15.12 | 0.1034 | 37.6723333 | -121.619 | Livermore | California | 2016-07-27 | 2016 | JUL | 27 | Wednesday | 00:19:43 |
| 30 | 4.9 | 371 | 97.07 | 1.439 | 21.5146 | 94.5721 | Pakokku | Burma | 2016-07-27 | 2016 | JUL | 27 | Wednesday | 00:20:28 |
| 122 | 0.4 | 2 | 1.09 | 0.02699 | 37.5958333 | -118.9948333 | Mammoth Lakes | California | 2016-07-27 | 2016 | JUL | 27 | Wednesday | 00:35:44 |
| 113.61 | 0.3 | 1 | 7.6 | 0.063 | 39.3775 | -119.845 | Mogul | Nevada | 2016-07-27 | 2016 | JUL | 27 | Wednesday | 00:41:59 |
| 50.399995968 | 1.8 | 50 | 1.3 | 0.04491576 | 61.2963 | -152.46 | Redoubt Volcano | Alaska | 2016-07-27 | 2016 | JUL | 27 | Wednesday | 00:52:52 |
| 93 | 1 | 15 | 2.452 | 0.0266 | 19.4235 | -155.6098333 | Honaunau-Napoopoo | Hawaii | 2016-07-27 | 2016 | JUL | 27 | Wednesday | 00:53:35 |
| 86.399993088 | 2 | 62 | 0.1 | 0.04581408 | 61.3019 | -152.4507 | Nikiski | Alaska | 2016-07-27 | 2016 | JUL | 27 | Wednesday | 00:58:45 |
| 149 | 1.2 | 22 | 0.18 | 0.1688 | 35.503 | -118.4058333 | Bodish | California | 2016-07-27 | 2016 | JUL | 27 | Wednesday | 01:03:00 |
| 87 | 1.67 | 43 | 12.91 | 0.106 | 37.673 | -121.6133333 | Livermore | California | 2016-07-27 | 2016 | JUL | 27 | Wednesday | 01:04:32 |

## Figure 2

*Month to Text Calculation*

```
1  case
2  when month_ = 7 then 'JUL'
3  when month_ = 8 then 'AUG'
4  end
```
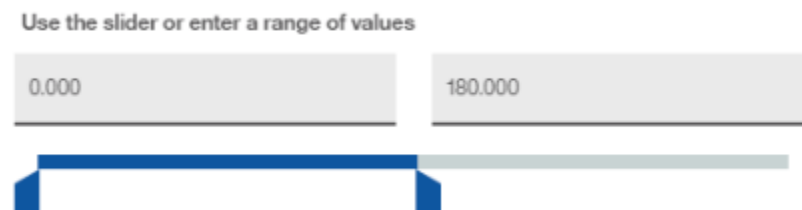
## Figure 3

*Gap Outlier filter*



## Figure 4
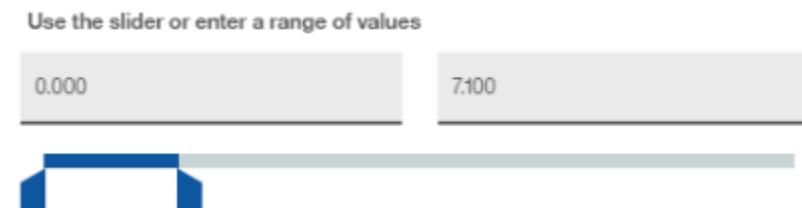
## Distance Outlier Filter

**Figure 5**

*Heat Map of Earthquake Occurrences*



**Figure 6**

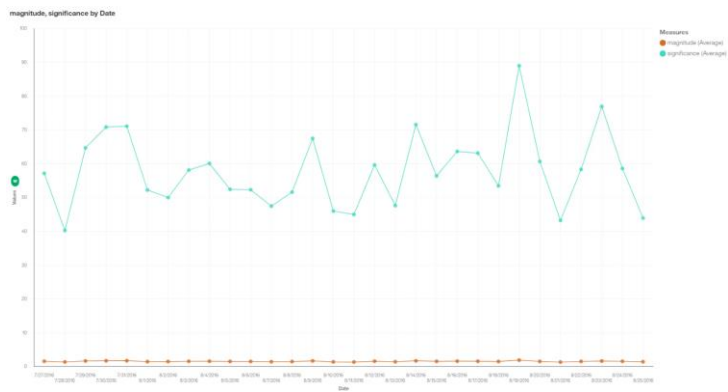*Magnitude and Significance by Date*



**Figure 7**
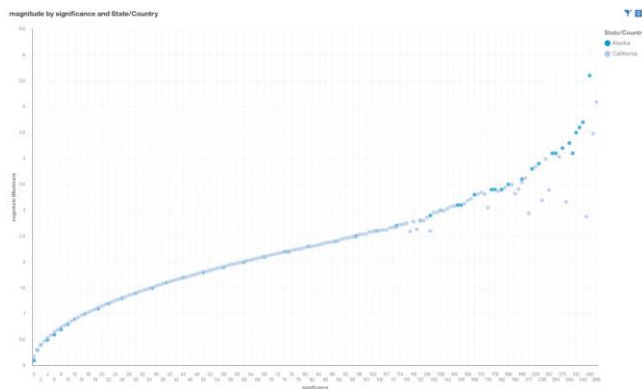
*Magnitude and significance of Alaska and California*
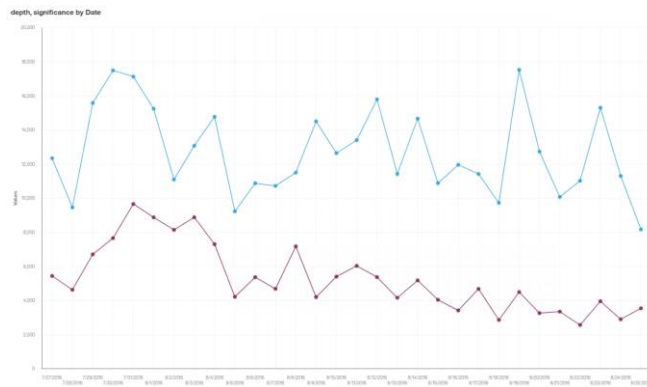
**Figure 8**

*Depth and Significance by Date*



**Figure 9**
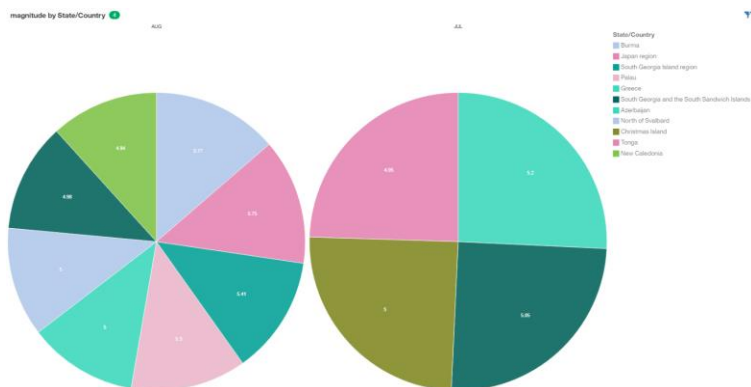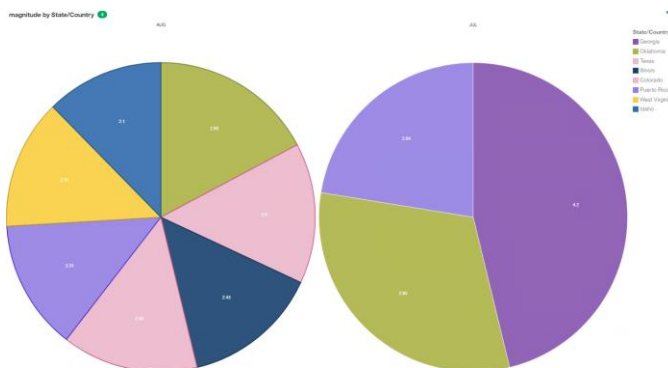
*Top 10 Magnitudes by County*



**Figure 10**

*Top 10 Magnitudes in the USA*

**Figure 11**

*Depth, Magnitude, and Significance by Month*



**Figure 12**

*Tree Diagram Using Magnitude as Target*



**Figure 13**

*Tree Diagram Using Distance as Target*

**Figure 14**

*Rules for Predicting Magnitude*



**Figure 15**

*Rules for Predicting Distance*



**Figure 16**

*Top 10 Earthquake locations by activity*
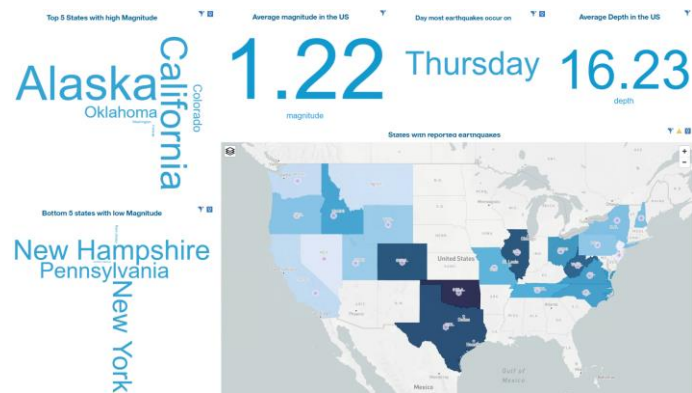
**Figure 17**

*US Earthquakes*



**Figure 18**

*World Earthquakes by Hemisphere*



**Figure 19**

*World Earthquake Location by Magnitude and Significance*