

#2 Group Project Report

IMAAN AGHA 1982051

RADHIKA SONDE 1964766

STEPHANY LOPEZ 2142613

ROBERTO JACKSON BAEZA 1861128

Summarization, Mapping, Hotspot Discovery and Change Analysis of High-Intensity Solar Flare Events

1 TASK 1

1.1 Method 1 - Flare intensity estimation based on the total.counts attribute

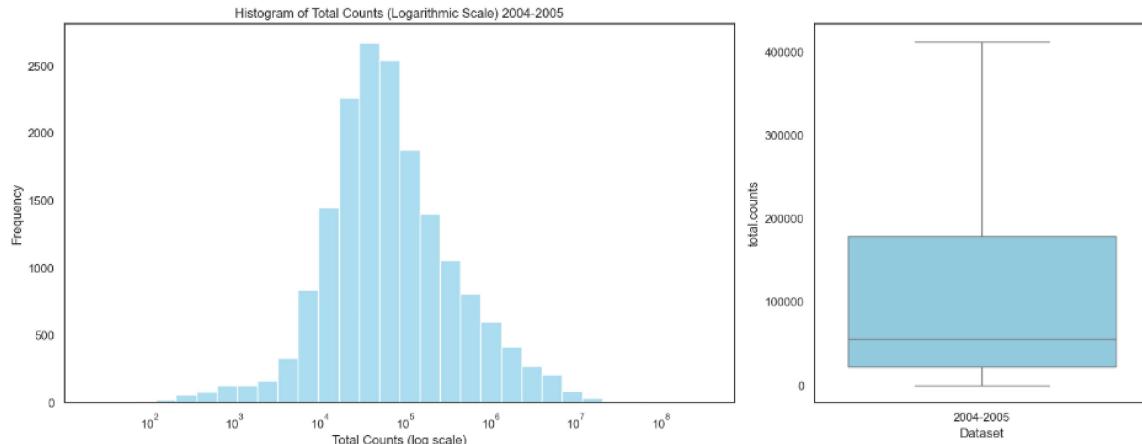


Fig. 1. Method 1 histogram and boxplot depicting total.counts

For Method 1, we struggled to visualize the total.counts data with so much noise from outliers. For our histogram, we found the best method to best understand the distribution of data was to use a logarithmic scale. When reading the histogram, however, it is important to realize that the increments on the x-axis are not uniform and increase exponentially. The Box plot was created by removing the outliers. We can clearly see the range of data without considering outliers. From these visualizations, we can tell that highest frequency of values is when the total.counts are in the range between 10^4 and 10^5 .

1.2 Method 2 - Flare intensity estimation based on the duration.s and energy.kev attributes

For Method 2, we created a bar graph to identify the number of occurrence of each energy band group. With the boxplot and the density plot, we were able to hone in on the relationship between energy bands and duration.

In 2004-2005, it seems that '6-12' energy bands have the highest frequency and the '300-800' energy bands have the lowest indicating that low energy flares occur more frequently than high energy flares. Conversely '300-800' energy bands have the longest duration while '6-12' have the shortest hinting at a relationship between both variables that as the energy levels increase, so does the duration of the flare.

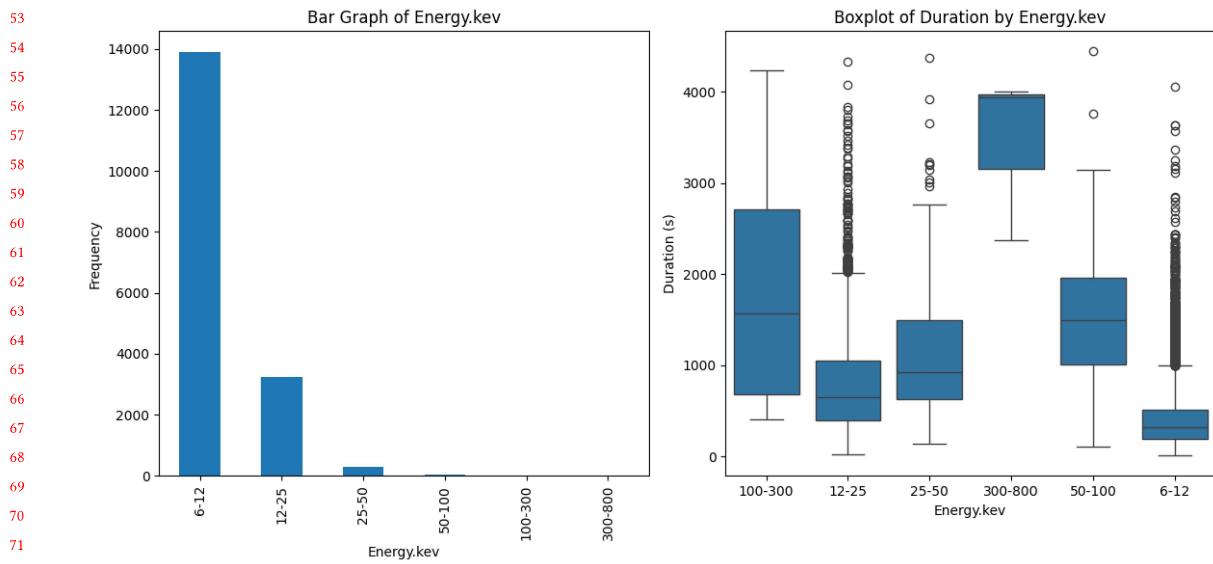


Fig. 2. Method 2 bar graph depicting the energy.kev frequency and boxplot illustrating the relationship between duration.s and energy.kev

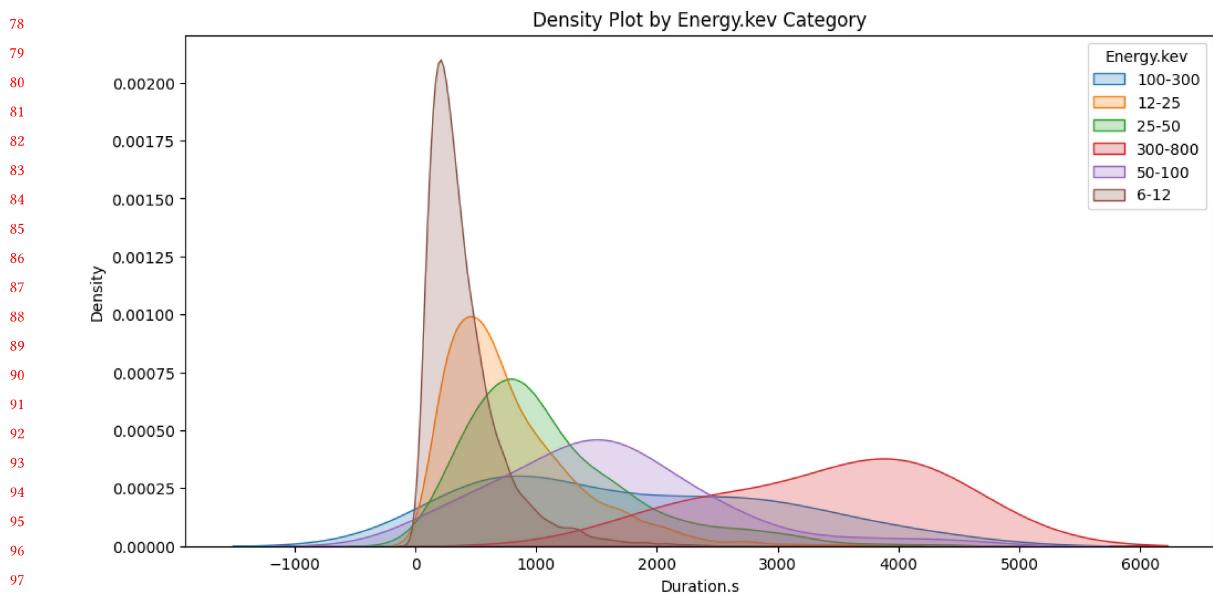


Fig. 3. Method 2 density depicting the relationship between duration and energy

1.3 Method 1 - Intensity heatmaps based on total.counts attribute

For Method 1, we generated two intensity heatmaps corresponding to months 1+2+3+4 and 21+22+23+24. In these heatmaps, the color of each pixel represents the intensity of solar flares at that position. The intensity is directly proportional to the values recorded in the 'total.counts' attribute.

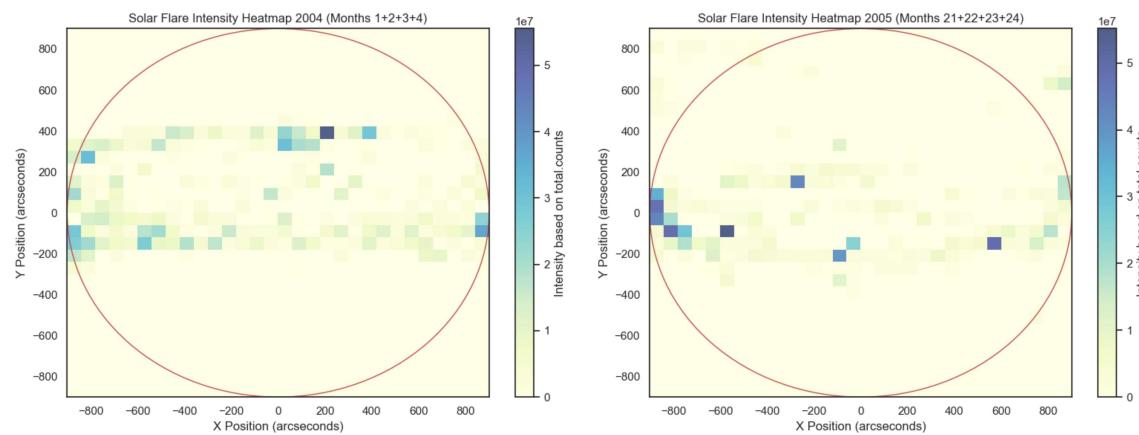


Fig. 4. Method 1 intensity heatmap based on total.counts

1.4 Method 2 - Intensity heatmaps based on duration.s and energy.kev attributes

For Method 2, the intensity estimation is computed using the formula:

$$\text{Intensity} = \text{duration.s} \times \frac{\text{energy.kev.f} + \text{energy.kev.i}}{2}$$

where the duration.s attribute represents the duration of the solar flare, and the energy.kev.i and energy.kev.f attributes contain the initial and final values of the energy range of the solar flare, respectively. The average energy of the solar flare is then multiplied by the duration to obtain the corresponding intensity. A logarithmic color scale had to be applied (LogNorm) to enhance the visualization of intensity difference. This was done because the intensity values mostly spanned in the low-intensity regions. The minimum and maximum values for the logarithmic scale (vmin and vmax) are calculated based on the minimum and maximum values of the intensity computed from the entire dataset. By incorporating a logarithmic scale, the resulting heatmaps effectively display the intensity variations.

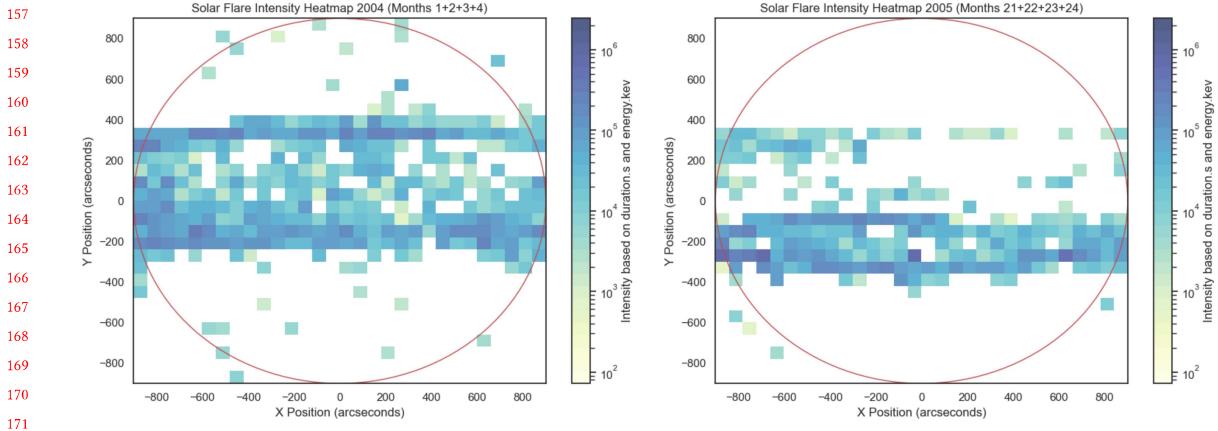


Fig. 5. Method 2 intensity heatmap based on duration.s and energy.kev

1.5 Intensity Map Comparison

Upon comparing the intensity maps based on total counts with those derived from duration and energy, we observe an overall higher solar flare intensity estimation for months 1+2+3+4 than for months 21+22+23+24. This suggests a correlation between the total counts, duration, and energy attributes. Additionally, we note that the intensity estimation of the solar flares is concentrated within position ranges of -400 to 400 on the Y-axis.

2 TASK 2

2.1 Hotspot Discovery Algorithm & Thresholds

```

187 from scipy.stats import gaussian_kde
188 import numpy as np
189 import matplotlib.pyplot as plt
190 import pandas as pd
191 import matplotlib.image as mpimg
192
193 # Read data file and set d1
194 df = pd.read_csv('Solar_flare_RHESSI_2004_05.csv')
195
196 # Create kernel density estimation
197 data = np.vstack([df['x.pos.asec'], df['y.pos.asec']])
198 kde = gaussian_kde(data)
199
200 # Create a gridspace and evaluate KDE
201 xgrid = np.linspace(-900, 900, 40)
202 ygrid = np.linspace(-900, 900, 40)
203 Xgrid, Ygrid = np.meshgrid(xgrid, ygrid)
204
205 grid_points = np.vstack([Xgrid.ravel(), Ygrid.ravel()])
206 Z = kde.evaluate(grid_points).reshape(Xgrid.shape)
207 thresholdZData = kde.evaluate(np.vstack([Xgrid.ravel(), Ygrid.ravel()]))
208 d1 = np.percentile(thresholdZData, 75) # Change this threshold value as needed
209 d2 = np.percentile(thresholdZData, 95) # Change this threshold value as needed
210

```

We developed the hotspot algorithm utilizing kernel density estimation to analyze the spatial distribution and precisely locate the hotspots on the sun using a 2D gridspace. After creating the regions, we filtered the data based on the selected threshold values. Subsequently, we generated a KDE plot and used the "coolwarm" color scheme to effectively pinpoint and enhance the visibility of these hotspots.

To determine our threshold values, d_1 and d_2 , we relied on the 75th and 95th percentiles of the density points. We opted for these percentiles because we recognize that hotspots consist of elevated concentrations of data points, surpassing the mean or 50th percentile. By selecting the 75th percentile, we effectively isolated the top 25% of our data, defining larger, more regionally significant hotspots with intensities exceeding a "medium-high" level. In contrast, the 95th percentile pinpointed the top 5% of our data, accurately representing smaller, highly concentrated hotspots with densities surpassing a "high" intensity threshold.

2.2 Hotspot Visualization

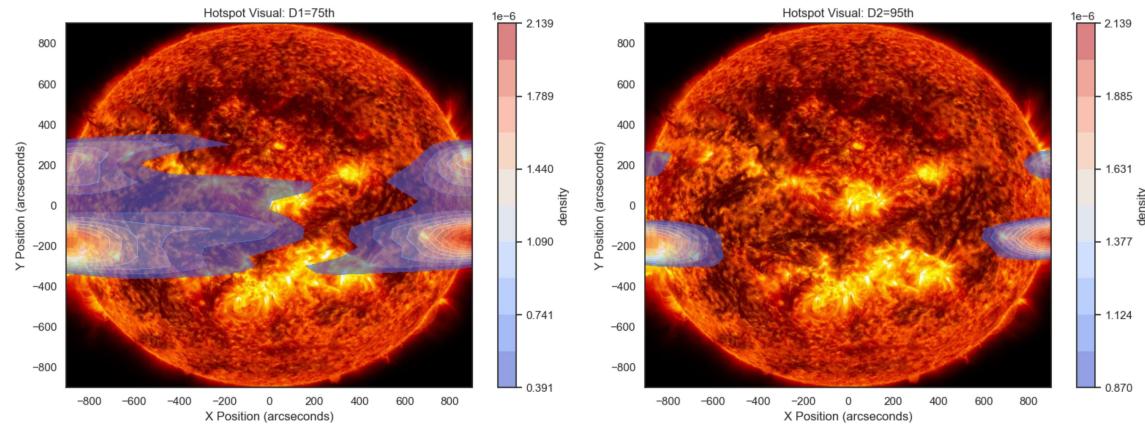
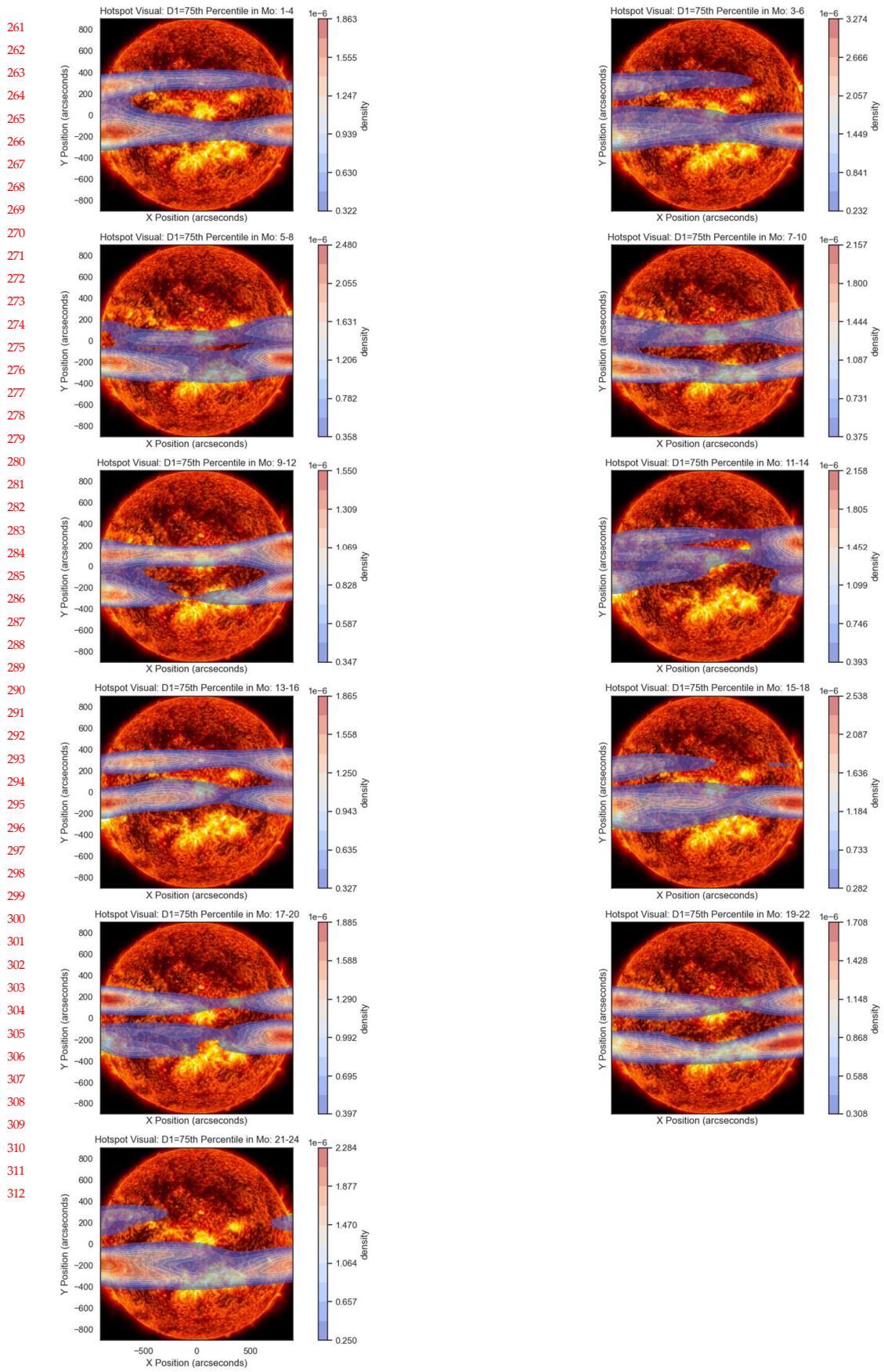
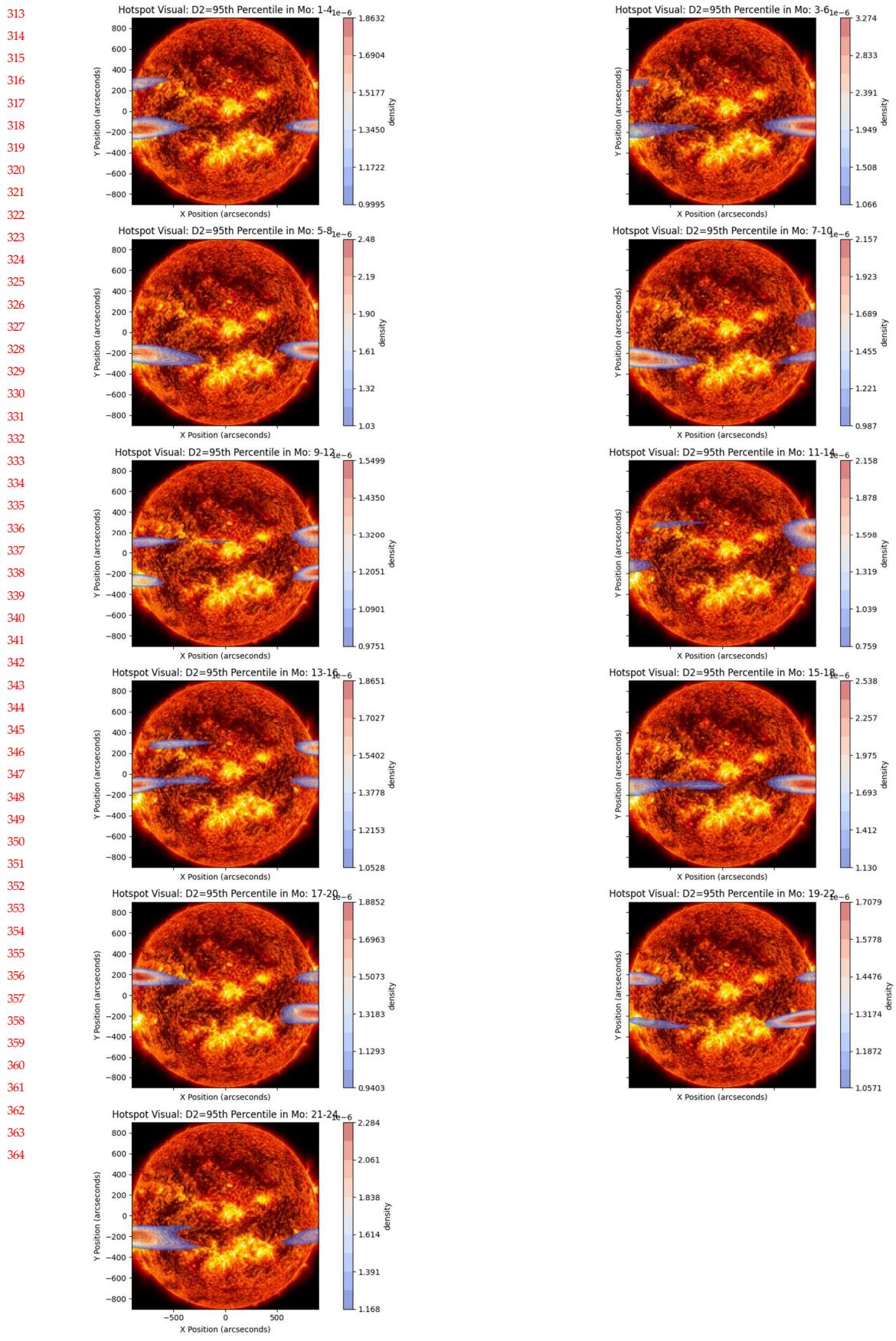


Fig. 6. Hotspot Visualization of the d_1 and d_2 threshold for all the months

2.3 Times Series Hotspot





365

366 To generate our time series, we employed our hotspot algorithm and filtered the data based on both the threshold
 367 value and batches of months.

368

369 In our time series for medium-high level, large hotspots, using a threshold value corresponding to the 75th percentile
 370 of the density points, we observe predominant hotspot activity both above and below the Sun's equator. Specifically, in
 371 the y-direction, activity spanned from -400 to 200, with a notable gap around the equator (-50 to 50) where no hotspots
 372 occurred. These hotspots align parallel to the x-axis, forming two distinct strips. Over time, while medium-high spots
 373 are consistently visible between 0 and -400 in the y direction, the activity above fluctuated. For example, in the initial
 374 batches, we observe low activity in the positive y-direction, which increases notably in the 9-12, 13-16, and 19-22 month
 375 batches.
 376

377

378 In the time series for smaller, highly intense hotspots, identified with a threshold value corresponding to the 95th
 379 percentile of the density points we observed that these hotspots primarily occurred at the far ends in the x-direction
 380 (<-600 and >600). Over time, similar to the medium-high hotspots, these also manifested between -400 to -50 and 50 to
 381 400 in the y-direction and in the southern region of the equator, with an exception in the 17-20 month batch where the
 382 hottest hotspot is found in the mid-northern region of the Sun.
 383

384

385 Based on this analysis, we can conclude that hotspots on the Sun predominantly occur in the mid-equator region as
 386 opposed to the poles, and for the hottest points, they additionally occur at the far corners of the Sun.

387

3 TASK 3

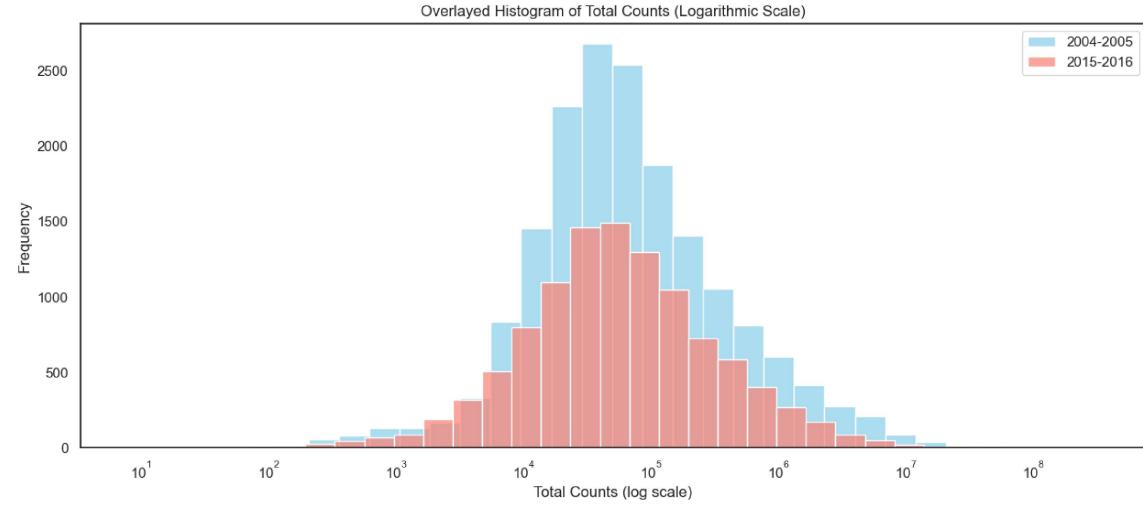
388

3.1 Change Analysis for Solar Flares Counts

389

390

391



392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

Fig. 7. Histogram comparison of total.counts between 2004-2005 & 2015-2016 dataset

412 It was hard to visualize the total.counts attribute due to the large number of outliers in the datasets. The total.counts
 413 attribute refers to the total number of counts in energy range 6-12keV over duration of a solar flare. The counts are
 414 corrected and summed over the entire duration of the flare event and including background noise. Using a logarithmic
 415

416

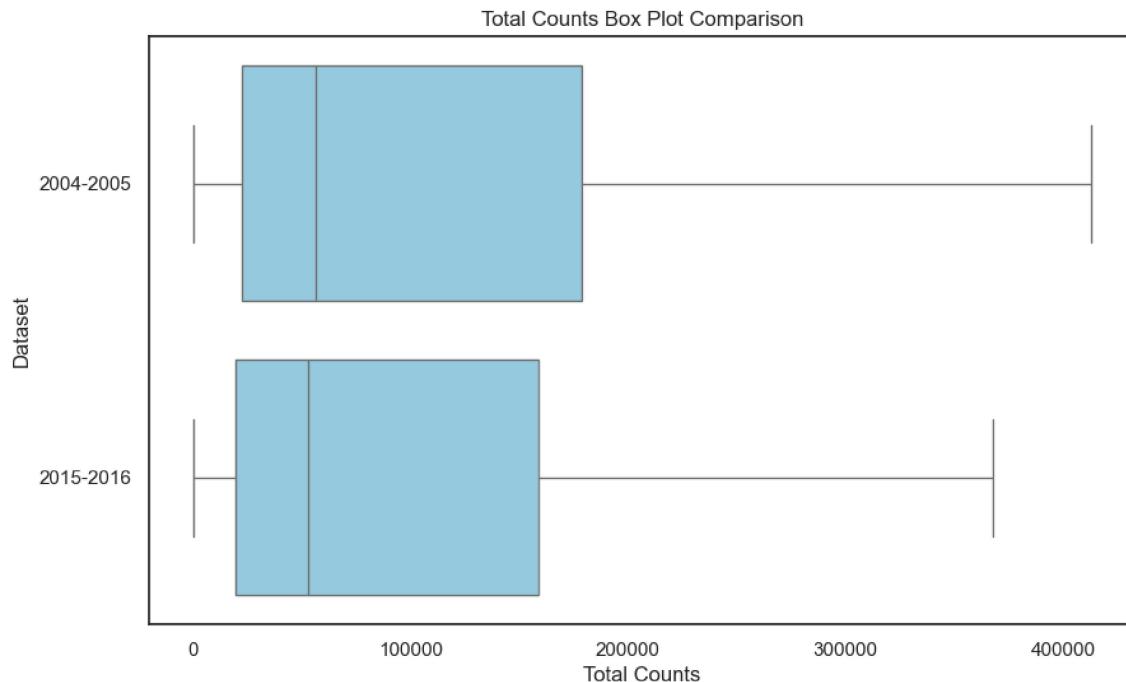


Fig. 8. Boxplot comparison of total.counts between 2004-2005 & 2015-2016 dataset

scale we were able to use a histogram to visualize the data. Using this technique helps to better understand the distribution of data when there are extreme or a wide range of values. The width of the bins increases exponentially and the distance between ticks are not equal intervals but rather ratios.

While the distribution of both histograms are similar the graph from 2004-05 has the higher frequency of events (over 2500). The histogram of 2015-16 has its peak frequency at just above 1400.

Plotting a boxplot was also a challenge with so many outliers, to make it easier to visualize the total.counts data we omitted the outliers to have reduced variability. Both boxplots have similar interquartile ranges but the 2004-05 plot has the bigger range of values that are not outliers.

3.2 Change Analysis for Solar Flares Duration and Energy Levels

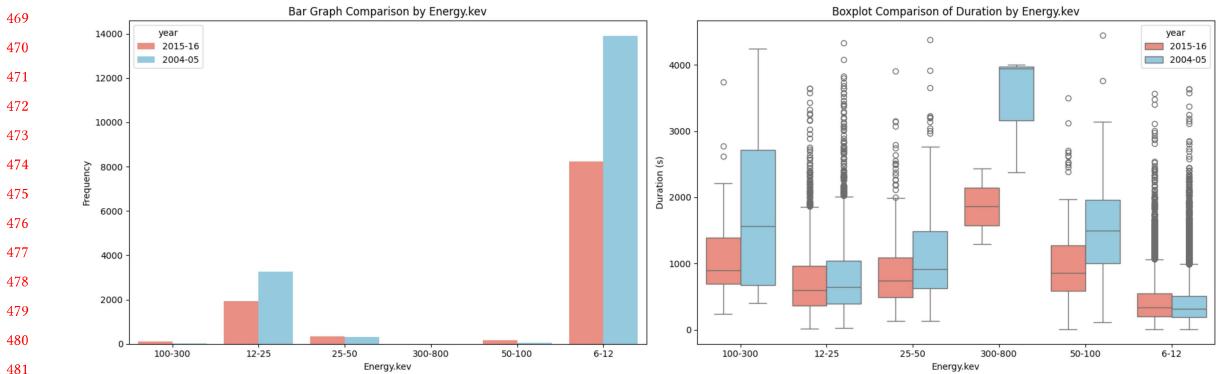


Fig. 9. Bar graph & Boxplot comparison of energy.kev and duration.s between 2004-2005 & 2015-2016 dataset

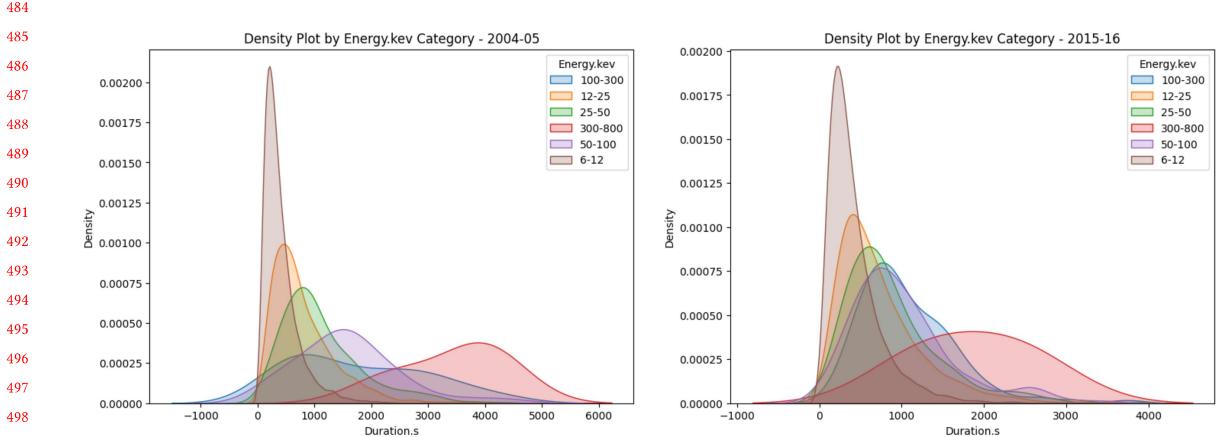


Fig. 10. Density plot comparison of energy.kev and duration.s between 2004-2005 & 2015-2016 dataset

Like the 2004-2005 dataset, the '6-12' energy bands have the highest frequency while the '300-800' energy bands have the lowest solidifying the observation that flares with high energy occur less frequently. Additionally, in the 2015-2016 period, the peaks of each energy category seem to be in order as the duration increases also bolstering the relationship between high energy levels having a longer duration than low ones.

From 2004-2005 to 2015-2016, there seems to be a major change in the energy levels and duration although the initial observations stay the same. The frequency of the energy bands decreased with the '6-12' band experiencing the greatest decline. There also seems to be a trend of decrease between the energy bands and the duration with '300-800', the highest energy level, experiencing the sharpest. Therefore, from 2004-2005 to 2015-2016, we can conclude the activity of the solar flares have decreased.

3.3 Change Analysis for Solar Flares Intensity

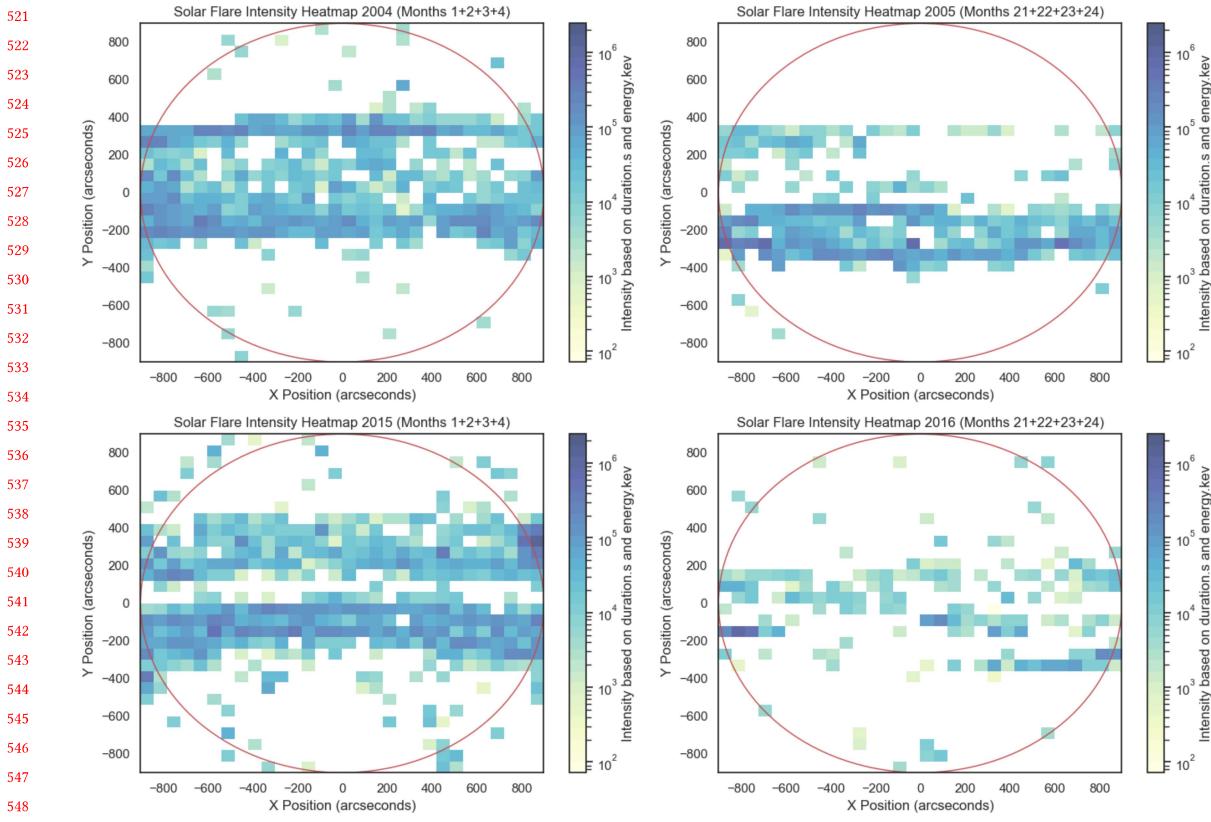


Fig. 11. Intensity heatmap comparison between 2004-2005 & 2015-2016 dataset

The solar flare intensity heatmaps from both datasets reveal higher intensities during the initial four months of 2004 and 2015 compared to the concluding four months of 2005 and 2016. Additionally, the distribution of solar flares is concentrated within the -400 to 400 range across all four plots. However, in the intensity maps for months 1+2+3+4 for 2004 and 2015, there are additional positions scattered outside the previously mentioned concentration range. The same phenomenon occurs for months 21+22+23+24 for the years 2005 and 2016 which could be an indicative of specific spatial patterns or anomalies in solar flare activity during those periods.

4 CONTRIBUTIONS

- **Imaan Agha** = Method 1 graphs for Task 1 & 3, Report
- **Roberto Jackson Baeza** = Method 1 intensity heatmaps, Times Series hotspots visualization & algorithm
- **Stephany Lopez** = Method 2 intensity heatmaps, intensity heatmaps comparison writeup, formatting time series graphs, Report
- **Radhika Sonde** = Method 2 graphs for Task 1 & 3, explanation summarization and interpretation of Task 2, Report

5 REFERENCES

- Intro to Spatial Analysis and HotSpot Discovery PowerPoint
- R Tutorial: Hotspot Analysis using Getis Ord Gi: https://rpubs.com/heatherleeleary/hotspot_getisOrd_tut
- COSC 3337 Group Project Introduction PowerPoint
- ChatGPT
- Exploratory Spatial Data Analysis (ESDA): https://darribas.org/gds_scipy16/ipynb_md/04_esda.html
- Exploratory Spatial and Temporal Data Analysis (ESTDA): https://darribas.org/gds_scipy16/ipynb_md/05_spatial_dynamics.html